

# Deep Dimension Reduction for Supervised Representation Learning

Jian Huang\*

Department of Statistics and Actuarial Science, University of Iowa, USA

Yuling Jiao†

School of Mathematics and Statistics, Wuhan University, China

Xu Liao

Center of Quantitative Medicine, Duke-NUS Medical School, Singapore

Jin Liu‡

Center of Quantitative Medicine, Duke-NUS Medical School, Singapore  
and

Zhou Yu

School of Statistics, East China Normal University, China

September 2, 2022

## Abstract

The goal of supervised representation learning is to construct effective data representations for prediction. Among all the characteristics of an ideal nonparametric representation of high-dimensional complex data, sufficiency, low dimensionality and disentanglement are some of the most essential ones. We propose a deep dimension reduction approach to learning representations with these characteristics. The proposed approach is a nonparametric generalization of the sufficient dimension reduction method. We formulate the ideal representation learning task as that of finding

---

\*Supported in part by the NSF grant DMS-1916199

†Supported in part by the NSFC grants No.11871474 and No.61701547

‡Supported in part by Duke-NUS Medical School WBS R-913-200-098-263 and Singapore MOE grants 2016-T2-2-029, 2018-T2-2-006 and 2018-T2-1-046

a nonparametric representation that minimizes an objective function characterizing conditional independence and promoting disentanglement at the population level. We then estimate the target representation at the sample level nonparametrically using deep neural networks. We show that the estimated deep nonparametric representation is consistent in the sense that its excess risk converges to zero. Our extensive numerical experiments using simulated and real benchmark data demonstrate that the proposed methods have better performance than several existing dimension reduction methods and the standard deep learning models in the context of classification and regression.

*Keywords:* Conditional independence; Distance covariance;  $f$ -divergence; Nonparametric estimation; Neural networks

# 1 Introduction

Over the past decade, deep learning has achieved impressive successes in modeling high-dimensional complex data arising from many scientific fields. A key factor for these successes is the ability of certain neural network models to learn nonlinear representations from complex high-dimensional data (Bengio et al., 2013; LeCun et al., 2015). For example, convolutional neural networks are able to learn effective representations of image data (LeCun et al., 1989). However, in general, optimizing the standard cross-entropy loss for classification and the least squares loss for regression do not guarantee that the learned representations enjoy any desired properties (Alain and Bengio, 2016). Therefore, it is imperative to develop principled approaches for constructing effective data representations. Representation learning has emerged as an important framework for modeling complex data (Bengio et al., 2013), with wide applications in classification, regression, imaging analysis, domain adaptation and transfer learning, among others. The goal of supervised representation learning is to construct effective representations of high-dimensional input data for various supervised learning tasks. In this paper, we propose a deep dimension reduction (DDR) method for sufficient representation learning. DDR aims at estimating a sufficient representation nonparametrically using deep neural networks based on the conditional independence principle.

There is a large body of literature on dimension reduction in statistics and machine learning. A prominent approach for supervised dimension reduction and representation learning is the sufficient dimension reduction (SDR) introduced in the seminal paper by Li (1991). A key aspect that distinguishes SDR from many other dimension reduction methods is that it does not make any model assumptions on the conditional distribution of the response given the predictors. In the framework of SDR, a semiparametric method,

called sliced inverse regression (SIR), was first proposed for estimating the linear dimension reduction direction, or linear sufficient representation (Li, 1991). The SIR and related methods were further developed by many researchers, see, for example, Cook and Weisberg (1991), Li (1992), Yin and Cook (2002), Cook (1998), Li et al. (2005) and Zhu et al. (2010) and the references therein. These methods require the linearity and constant covariance conditions on the distribution of the predictors. Several approaches have been developed without assuming these conditions, including methods based on nonparametric regression (Xia et al., 2002), conditional covariance operators (Fukumizu et al., 2009), mutual information (Suzuki and Sugiyama, 2013), distance correlation (Vepakomma et al., 2018), and semiparametric modeling (Ma and Zhu, 2012, 2013a). These SDR methods focus on linear dimension reduction, that is, the features learned are linear functions of the original input variables. However, linear functions may not be adequate for representing high-dimensional complex data such as images and natural languages, due to the highly nonlinear nature of such data. Lee et al. (2013) formulated a general sufficient dimension reduction framework in the nonlinear setting and proposed a generalized inverse regression approach using conditional covariance operators, but this method is computationally prohibitive with high-dimensional data such as the image datasets considered in Section 6. We refer to the review papers (Cook, 2007, 2018; Ma and Zhu, 2013b) and the monograph (Li, 2018) for thorough reviews of SDR methods.

Among all the characteristics of an ideal representation for supervised learning, sufficiency, low dimensionality and disentanglement are some of the most essential ones (Achille and Soatto, 2018). Sufficiency is a basic property a representation should have. It is closely related to the concept of sufficient statistics in a parametric model (Fisher, 1922; Cook, 2007). In supervised representation learning, sufficiency is characterized by the conditional independence principle, which states that the original input data is conditionally indepen-

dent of the response given the representation. In other words, a sufficient representation contains all the relevant information in the input data about the response. Low dimensionality means that the representation should have as few components as possible to represent the underlying structure of the data, and the number of components should be fewer than the ambient dimension. In the context of nonparametric representation learning, disentanglement refers to the requirement that the components of the representation should be statistically independent. This is an extension of and stronger than the orthogonal constraint in the linear representation setting, where the components of the linear representation are constrained to have orthonormal directions. The notion of disentanglement is based on the hypothesis that there are some underlying factors determining the data generation process: although the observed data are high-dimensional and complex, the underlying factors are low-dimensional, disentangled, and have a simple statistical structure. The components in the learned representation can often be interpreted as corresponding to the latent structure of the observed data, thus disentanglement is an important property for better separating latent factors from one to another. A representation with these characteristics can make the model more interpretable and facilitates the downstream supervised learning tasks.

Inspired by the basic idea of SDR, we propose a deep dimension reduction (DDR) approach for supervised representation learning with the properties of sufficiency, low dimensionality and disentanglement. By taking the advantage of the strong capacities of deep neural networks in approximating high-dimensional functions for nonparametric estimation, we model the DDR representations, which we refer to as DDR map (DDRM) for convenience, using deep neural networks to capture the nonlinearity in the representation space. It would be difficult to use the traditional techniques for nonparametric estimation such as kernel smoothing and splines for multi- or high-dimensional function estimation in the context of representation learning. To characterize the conditional independence of

the representation, we use the distance covariance (Székely et al., 2007) as the conditional independence measure that can be computed efficiently. We also promote the disentanglement for DDRM by regularizing its distribution to have independent components based on a divergence measure.

Our main contributions are as follows:

- We formulate a new nonparametric approach to dimension reduction by characterizing the sufficient dimension reduction map as a minimizer of a loss function measuring conditional independence and disentanglement.
- We estimate the sufficient dimension reduction map at the sample level nonparametrically using deep neural networks based on distance covariance for characterizing sufficiency and use  $f$ -divergence to promote disentanglement of the learned representation.
- We show that the estimated deep dimension reduction map is consistent in the sense that it achieves asymptotic sufficiency under mild conditions.
- We validate DDR via comprehensive numerical experiments and real data analysis in the context of regression and classification. We use the learned features based on DDR as inputs for linear regression and nearest neighbor classification. The resulting prediction accuracies are better than those based on linear dimension reduction methods for regression and deep learning models for classification. The PyTorch code for DDR is available at <https://github.com/anonymous/DDR>.

The rest of the paper is organized as follows. In Section 2 we discuss the theoretical framework for learning a DDRM. This framework leads to the formulation of an objective function using distance correlation for characterizing conditional independence in Section 3. We estimate the target DDRM based on the sample version of the objective function using

deep neural networks and develop an efficient algorithm for training the DDRM. In Section 4 we provide sufficient conditions under which estimated nonparametric representations achieves asymptotic sufficiency. This result provides strong theoretical support for the proposed method. The algorithm for implementing DDR is described in Section 5. In Section 6 we validate the proposed DDR via extensive numerical experiments and real data examples.

## 2 Sufficient representation and distance correlation

Consider a pair of random vectors  $(X, Y) \in \mathbb{R}^p \times \mathbb{R}^q$ , where  $X$  is a vector of predictors and  $Y$  is a vector of response variables or labels. Our goal is to construct a representation of  $X$  that possesses the three characteristics: sufficiency, low dimensionality and disentanglement.

### 2.1 Sufficiency

A measurable function  $\mathbf{s} : \mathbb{R}^p \rightarrow \mathbb{R}^d$  with  $d \leq p$  is said to be a sufficient representation of  $X$  if

$$Y \perp\!\!\!\perp X | \mathbf{s}(X), \tag{1}$$

that is,  $Y$  and  $X$  are conditionally independent given  $\mathbf{s}(X)$ . This condition holds if and only if the conditional distribution of  $Y$  given  $X$  and that of  $Y$  given  $\mathbf{s}(X)$  are equal. Therefore, the information in  $X$  about  $Y$  is completely encoded by  $\mathbf{s}(X)$ . Such a function  $\mathbf{s}$  always exists, since if we simply take  $\mathbf{s}(\mathbf{x}) = \mathbf{x}$ , then (1) holds trivially. This formulation is a nonparametric generalization of the basic condition in sufficient dimension reduction (Li, 1991; Cook, 1998), where it is assumed  $\mathbf{s}(\mathbf{x}) = \mathbf{B}^T \mathbf{x}$  with  $\mathbf{B} \in \mathbb{R}^{p \times d}$  belonging to the Stiefel manifold, i.e.,  $\mathbf{B}^T \mathbf{B} = \mathbf{I}_d$ .

Denote the class of sufficient representations satisfying (1) by

$$\mathcal{F} = \{\mathbf{s} : \mathbb{R}^p \rightarrow \mathbb{R}^d, \mathbf{s} \text{ satisfies } Y \perp\!\!\!\perp X | \mathbf{s}(X)\}.$$

For an injective measurable transformation  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$  and  $\mathbf{s} \in \mathcal{F}$ ,  $T \circ \mathbf{s}(X)$  is also sufficient by the basic property of conditional probability. Therefore, the class  $\mathcal{F}$  is invariant in the sense that

$$T \circ \mathcal{F} \subseteq \mathcal{F}, \text{ provided } T \text{ is injective,}$$

where  $T \circ \mathcal{F} = \{T \circ \mathbf{s} : \mathbf{s} \in \mathcal{F}\}$ . An important class of transformations is the class of affine transformations,  $T \circ \mathbf{s} = \mathbf{A}\mathbf{s} + \mathbf{b}$ , where  $\mathbf{A}$  is a  $d \times d$  nonsingular matrix and  $\mathbf{b} \in \mathbb{R}^d$ .

## 2.2 Space of nonparametric sufficient representations

The nonparametric sufficient representations are nonunique and the space of such representations is large, since if (1) holds for  $\mathbf{s}$ , it also holds for any one-to-one transformation of  $\mathbf{s}$ . We propose to narrow the space of such representations by constraining the distributional properties of  $\mathbf{s}(\mathbf{x})$ .

Among the sufficient representations, it is preferable to have those with a simple statistical distribution and whose components are independent, that is, the components are disentangled. For a sufficient representation  $\mathbf{s}(X)$ , let  $\Sigma_{\mathbf{s}} = \text{Cov}(\mathbf{s}(X))$ . Suppose  $\Sigma_{\mathbf{s}}$  is positive definite, then  $\Sigma_{\mathbf{s}}^{-1/2}\mathbf{s}(X)$  is also a sufficient representation. Therefore, we can always rescale  $\mathbf{s}(X)$  such that it has identity covariance matrix. To further simplify the statistical structure of a representation  $\mathbf{s}$ , we also impose the constraint that it is rotation invariant in distribution, that is,  $\mathbf{Q}\mathbf{s}(X) = \mathbf{s}(X)$  in distribution for any orthogonal matrix  $\mathbf{Q} \in \mathbb{R}^{d \times d}$ . By the Maxwell characterization of the Gaussian distributions (Maxwell, 1860; Bryc, 1995),

a random vector of dimension two or more with independent components is rotation invariant in distribution if and only if it is Gaussian with zero mean and a spherical covariance matrix. Therefore, after absorbing the scaling factor, for a sufficient representation map to have independent components and be rotation invariant, it is necessarily distributed as  $N_d(\mathbf{0}, \mathbf{I}_d)$ . Denote

$$\mathcal{M} = \{R : \mathbb{R}^p \rightarrow \mathbb{R}^d, R(X) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)\}. \quad (2)$$

Now our problem becomes that of finding a representation in  $\mathcal{F} \cap \mathcal{M}$ , the intersection of the Fisher class and the Maxwell class.

Does such a sufficient representation exist? The following result from the optimal transport theory gives an affirmative answer and guarantees the existence of such a representation under mild conditions (Villani, 2008).

**Lemma 2.1.** *Let  $\mu$  be a probability measure on  $\mathbb{R}^d$ . Suppose it has finite second moment and is absolutely continuous with respect to the standard Gaussian measure, denoted by  $\gamma_d$ . Then it admits a unique optimal transportation map  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that  $T_{\#}\mu = \gamma_d \equiv \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , where  $T_{\#}\mu$  denotes the pushforward distribution of  $\mu$  under  $T$ . Moreover,  $T$  is injective  $\mu$ -almost everywhere.*

Denote the law of a random vector  $Z$  by  $\mu_Z$ . Lemma 2.1 implies that, for any  $\mathbf{s} \in \mathcal{F}$  with  $\mathbb{E}\|\mathbf{s}(X)\|^2 < \infty$  and  $\mu_{\mathbf{s}(X)}$  absolutely continuous with respect to  $\gamma_d$ , there exists a map  $T^*$  transforming the distribution of  $\mathbf{s}(X)$  to  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Therefore,  $R^* := T^* \circ \mathbf{s} \in \mathcal{F} \cap \mathcal{M}$ , that is,

$$X \perp\!\!\!\perp Y | R^*(X) \quad \text{and} \quad R^*(X) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d). \quad (3)$$

The requirement that  $R^*(X) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  can be considered a regularization on the distribution of  $R^*(X)$ . This is similar to the ridge regression where the ridge penalty can be derived from a spherical normal prior on the regression coefficient. We use the standard

multivariate normal distribution as the reference for regularizing the distribution of the sufficient representation. It is possible to use other distributions such as uniform distribution on the unit cube  $[0, 1]^d$ . Below, to be specific, we will focus on using the standard normal distribution for regularization. Also, we note that it suffices to estimate the function  $R^*$ , not  $\mathbf{s}$  and  $T^*$  separately, since  $R^*$  satisfies the conditional independence requirement.

The independence requirement for the components of  $R^*$  is reminiscent of the same requirement in the independent component analysis (ICA, Jutten and Herault (1991); Comon (1994)). ICA is a method for estimating hidden factors that underlie a random vector  $X$ . It posits that  $X$  is a *linear transformation* of an unknown random vector with independent components, but the transformation is unknown. The goal of ICA is to estimate this linear transformation. DDR differs from ICA in three crucial aspects. First, DDR is a supervised method that seeks to find a data representation such that the response is conditionally independent given this representation, while ICA is an unsupervised method that attempts to identify independent latent factors underlying the original data vector. Second, DDR seeks a nonparametric function  $R^*$  such that  $R^*(X)$  has independent components, while ICA attempts to find a matrix  $W \in \mathbb{R}^{p \times p}$  such that the components of  $WX$  are independent. Third, the distribution of  $R^*(X)$  can be Gaussian; in contrast, a basic restriction in ICA is that the independent components must be non-Gaussian. There is a large body of literature on ICA. For some more recent references on ICA, see Samarov and Tsybakov (2004), Samworth and Yuan (2012) and the review Nordhausen and Oja (2018). Feedforward neural networks and recurrent neural network structures have also been considered in solving ICA problems (Mutihac and Hulle, 2003). We refer the reader to the monographs (Aapo Hyvärinen et al., 2001; Roberts and Everson, 2001) for additional references on ICA.

### 3 Nonparametric estimation of representation map

The discussions in Section 2 lay the ground for formulating an objective function that can be used for constructing a DDRM  $R^*$  satisfying (3), that is,  $R^*$  is sufficient and disentangled.

#### 3.1 Population objective function

Let  $\mathcal{V}$  be a measure of dependence between random variables  $X$  and  $Y$  with the following properties: (a)  $\mathcal{V}[X, Y] \geq 0$  with  $\mathcal{V}[X, Y] = 0$  if and only if  $X \perp\!\!\!\perp Y$ ; (b)  $\mathcal{V}[X, Y] \geq \mathcal{V}[R(X), Y]$  for all measurable function  $R$ ; and (c)  $\mathcal{V}[X, Y] = \mathcal{V}[R^*(X), Y]$  if and only if  $R^* \in \mathcal{F}$ . These properties imply that  $R^* \in \mathcal{F}$  if and only if  $R^* \in \operatorname{argmin}_R \{-\mathcal{V}[R(X), Y]\}$ .

For the normality regularization in (3), we use a divergence measure  $\mathbb{D}$  to quantify the difference between  $\mu_{R(X)}$  and the standard normal distribution  $\gamma_d$ . This measure should satisfy the condition  $\mathbb{D}(\mu_{R(X)} \parallel \gamma_d) \geq 0$  for every measurable function  $R$  and  $\mathbb{D}(\mu_{R(X)} \parallel \gamma_d) = 0$  if and only if  $R \in \mathcal{M}$ . The  $f$ -divergences, including the KL-divergence, satisfy this condition. It follows that  $R^* \in \mathcal{M}$  if and only if  $R^* \in \operatorname{argmin}_R \mathbb{D}(\mu_{R(X)} \parallel \gamma_d)$ . Then the problem of finding a sufficient and disentangle map  $R^*$  becomes a constrained minimization problem:

$$\operatorname{argmin}_R -\mathcal{V}[R(X), Y] \quad \text{subject to} \quad \mathbb{D}(\mu_{R(X)} \parallel \gamma_d) = 0.$$

The Lagrangian form of this minimization problem is

$$\mathcal{L}(R) = -\mathcal{V}[R(X), Y] + \lambda \mathbb{D}(\mu_{R(X)} \parallel \gamma_d), \tag{4}$$

where  $\lambda \geq 0$  is a tuning parameter. This parameter provides a balance between the sufficiency property and the disentanglement constraint. A small  $\lambda$  leads to a representation with more emphasis on sufficiency, while a large  $\lambda$  yields a representation with more em-

phasis on disentanglement. We show in Theorem 3.2 below that any  $R^*$  satisfying (3) is a minimizer of  $\mathcal{L}(R)$ . Therefore, we can train a DDRM by minimizing an empirical version of  $\mathcal{L}(R)$ .

There are several options for  $\mathcal{V}$  with the properties (a)-(c) described above. For example, we can take  $\mathcal{V}$  to be the mutual information. However, in addition to estimating the DDRM  $R$ , this choice requires nonparametric estimation of the ratio of the joint density and the marginal densities of  $Y$  and  $R(X)$ , which is not an easy task. To be specific, in this work we use the distance covariance (Székely et al., 2007) between  $Y$  and  $R(X)$ , which has an elegant  $U$ -statistic expression. It does not involve additional unknown quantities and is easy to compute. For the divergence measure of two distributions, we use the  $f$ -divergence (Ali and Silvey, 1966), which includes the KL-divergence as a special case.

## 3.2 Empirical objective function

In this subsection, we formulate the objective function for the proposed deep dimension reduction method. We first describe some essentials about distance covariance and  $f$ -divergence.

### 3.2.1 Distance covariance

We recall the concept of distance covariance (Székely et al., 2007), which characterizes the dependence of two random variables. Let  $\mathbf{i}$  be the imaginary unit  $(-1)^{1/2}$ . For any  $\mathbf{t} \in \mathbb{R}^d$  and  $\mathbf{s} \in \mathbb{R}^m$ , let  $\psi_Z(\mathbf{t}) = \mathbb{E}[\exp \mathbf{i} \mathbf{t}^T Z]$ ,  $\psi_Y(\mathbf{s}) = \mathbb{E}[\exp \mathbf{i} \mathbf{s}^T Y]$ , and  $\psi_{Z,Y}(\mathbf{t}, \mathbf{s}) = \mathbb{E}[\exp \mathbf{i}(\mathbf{t}^T Z + \mathbf{s}^T Y)]$  be the characteristic functions of random vectors  $Z \in \mathbb{R}^d, Y \in \mathbb{R}^m$ , and the pair  $(Z, Y)$ ,

respectively. The squared distance covariance  $\mathcal{V}[Z, Y]$  is defined as

$$\mathcal{V}[Z, Y] = \int_{\mathbb{R}^{d+m}} \frac{|\psi_{Z,Y}(\mathbf{t}, \mathbf{s}) - \psi_Z(\mathbf{t})\psi_Y(\mathbf{s})|^2}{c_d c_m \|\mathbf{t}\|^{d+1} \|\mathbf{s}\|^{q+1}} d\mathbf{t} d\mathbf{s},$$

where  $c_d = \frac{\pi^{(d+1)/2}}{\Gamma((d+1)/2)}$ . Given  $n$  i.i.d copies  $\{Z_i, Y_i\}_{i=1}^n$  of  $(Z, Y)$ , an unbiased estimator of  $\mathcal{V}$  is the empirical distance covariance  $\widehat{\mathcal{V}}_n$ , which can be elegantly expressed as a  $U$ -statistic (Huo and Székely, 2016)

$$\widehat{\mathcal{V}}_n[Z, Y] = \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} h((Z_{i_1}, Y_{i_1}), \dots, (Z_{i_4}, Y_{i_4})),$$

where  $h$  is the kernel defined by

$$\begin{aligned} h((\mathbf{z}_1, \mathbf{y}_1), \dots, (\mathbf{z}_4, \mathbf{y}_4)) &= \frac{1}{4} \sum_{\substack{1 \leq i, j \leq 4 \\ i \neq j}} \|\mathbf{z}_i - \mathbf{z}_j\| \|\mathbf{y}_i - \mathbf{y}_j\| + \frac{1}{24} \sum_{\substack{1 \leq i, j \leq 4 \\ i \neq j}} \|\mathbf{z}_i - \mathbf{z}_j\| \sum_{\substack{1 \leq i, j \leq 4 \\ i \neq j}} \|\mathbf{y}_i - \mathbf{y}_j\| \\ &\quad - \frac{1}{4} \sum_{i=1}^4 \left( \sum_{\substack{1 \leq j \leq 4 \\ j \neq i}} \|\mathbf{z}_i - \mathbf{z}_j\| \sum_{\substack{1 \leq j \leq 4 \\ i \neq j}} \|\mathbf{y}_i - \mathbf{y}_j\| \right). \end{aligned}$$

For a categorical response  $Y$  in multi-class classification problems, we can use one-hot vectors to code the classes, i.e., for the  $k$ th class,  $Y$  is a unit vector with  $k$ th element equaling 1 and the remaining elements being 0. The  $L_2$  distance between two observed responses  $y_i$  and  $y_j$  is

$$\|y_i - y_j\|_2 = \begin{cases} 0, & \text{if } y_i = y_j, \\ \sqrt{2}, & \text{if } y_i \neq y_j. \end{cases}$$

Note that the number  $\sqrt{2}$  simply scales the whole objective function and does not affect the solution.

### 3.2.2 $f$ -divergence

Let  $\mu$  and  $\gamma$  be two probability measures on  $\mathbb{R}^d$ . The  $f$ -divergence (Ali and Silvey, 1966) between  $\mu$  and  $\gamma$  with  $\mu \ll \gamma$  is defined as

$$\mathbb{D}_f(\mu||\gamma) = \int_{\mathbb{R}^d} f\left(\frac{d\mu}{d\gamma}\right) d\gamma, \quad (5)$$

where  $f : \mathbb{R}^+ \rightarrow \mathbb{R}$  is a differentiable convex function satisfying  $f(1) = 0$ . Let  $f^*$  be the Fenchel conjugate of  $f$  (Rockafellar, 1970), defined by

$$f^*(t) = \sup_{x \in \mathbb{R}} \{tx - f(x)\}, t \in \mathbb{R}. \quad (6)$$

The  $f$ -divergence (5) admits the following variational formulation (Keziou, 2003; Nguyen et al., 2010; Nowozin et al., 2016).

**Lemma 3.1.** *Suppose that  $f$  is a differentiable convex function. Then,*

$$\mathbb{D}_f(\mu||\gamma) = \max_{D: \mathbb{R}^d \rightarrow \text{dom}(f^*)} \mathbb{E}_{Z \sim \mu} D(Z) - \mathbb{E}_{W \sim \gamma} f^*(D(W)), \quad (7)$$

where  $f^*$  is defined in (6). In addition, the maximum is attained at  $D(\mathbf{z}) = f'\left(\frac{d\mu}{d\gamma}(\mathbf{z})\right)$ .

Commonly used divergence measures include the Kullback-Leibler (KL) divergence, the Jensen-Shanon (JS) divergence and the  $\chi^2$ -divergence. We summarize the details in Table 1.

To be specific, in this paper we use the KL divergence with  $f(x) = x \log x$ , which has the familiar form  $\mathbb{D}_{\text{KL}}(\mu||\gamma) = \int_{\mathbb{R}^d} \left(\log \frac{d\mu}{d\gamma}\right) d\mu$ . The dual form of  $f$  is  $f^*(t) = \exp(t - 1)$ . The variational representation  $\mathbb{D}_{\text{KL}}(\mu||\gamma) = \sup_D \{\mathbb{E}_{Z \sim \mu} D(Z) - \mathbb{E}_{W \sim \gamma} \exp(D(W) - 1)\}$ . The generative adversarial networks (GAN, Goodfellow et al. (2014)) corresponds to the JS-divergence. Much work has been devoted to developing various extensions and alternative

Table 1: Three examples of  $f$ -divergence

$f$ -Div	$f(x)$	$f^*(t)$	$\mathbb{D}_f(\mu, \gamma)$
KL	$x \log x$	$e^{t-1}$	$\sup_D \{\mathbb{E}_{Z \sim \mu} D(Z) - \mathbb{E}_{W \sim \gamma} e^{D(W)-1}\}$
JS	$-(x+1) \log \frac{x+1}{2} + x \log x$	$-\log(2 - \exp(t))$	$\sup_D \{\mathbb{E}_{Z \sim \mu} D(Z) + \mathbb{E}_{W \sim \gamma} \log(2 - \exp(D(W)))\}$
$\chi^2$	$(x-1)^2$	$t + \frac{t^2}{4}$	$\sup_D \{\mathbb{E}_{Z \sim \mu} D(Z) - \mathbb{E}_{W \sim \gamma} [D(W) + \frac{D^2(W)}{4}]\}$

formulations of the original GAN (Li et al., 2015; Nowozin et al., 2016; Sutherland et al., 2017; Arjovsky et al., 2017).

### 3.2.3 Empirical objective function for DDR

We are now ready to formulate an empirical objective function for learning DDRM. Let  $R \in \mathcal{M}$ , where  $\mathcal{M}$  is the Maxwell class defined in (2). By the variational formulation (7), we can write the population version of the objective function (4) as

$$\mathcal{L}(R) = -\mathcal{V}[R(X), Y] + \lambda \max_D \{\mathbb{E}_{X \sim \mu_X} D(R(X)) - \mathbb{E}_{W \sim \gamma_d} f^*(D(W))\}. \quad (8)$$

This expression is convenient since we can simply replace the expectations by the corresponding empirical averages.

**Theorem 3.2.** *We have  $R^* \in \arg \min_{R \in \mathcal{M}} \mathcal{L}(R)$  provided (3) holds.*

According to Theorem 3.2, it is natural to estimate  $R^*$  based on the empirical version of the objective function (8) when a random sample  $\{(X_i, Y_i)\}_{i=1}^n$  is available.

We estimate  $R^*$  nonparametrically using feedforward neural networks (FNN) (Schmidhuber, 2015). Two networks are employed: the representer network  $R_\theta$  with parameter  $\theta$  for estimating  $R^*$  and a second network  $D_\phi$  with parameter  $\phi$  for estimating the discrim-

inator  $D$ . For any function  $f(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}^d$ , denote  $\|f\|_\infty = \sup_{\mathbf{x} \in \mathcal{X}} \|f(\mathbf{x})\|$ , where  $\|\cdot\|$  is the Euclidean norm.

- **Representer network  $R_\theta$** : This network is used for training  $R^*$ . Let  $\mathbf{R} \equiv \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  be the set of such ReLU neural networks  $R_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^d$  with parameter  $\theta$ , depth  $\mathcal{H}$ , width  $\mathcal{W}$ , size  $\mathcal{S}$ . Here the depth  $\mathcal{H}$  refers to the number of hidden layers, so the network has  $\mathcal{H} + 1$  layers in total. A  $(\mathcal{H} + 1)$ -vector  $(w_0, w_1, \dots, w_{\mathcal{H}})$  specifies the width of each layer, where  $w_0 = p$  is the dimension of the input data and  $w_{\mathcal{H}} = d$  is the dimension of the output. The width  $\mathcal{W} = \max\{w_1, \dots, w_{\mathcal{H}}\}$  is the maximum width of the hidden layers. The size  $\mathcal{S} = \sum_{i=0}^{\mathcal{H}} [w_i \times (w_i + 1)]$  is the total number of parameters in the network.
- **Discriminator network  $D_\phi$** : This network is used as the witness function for checking whether the distribution of the estimator of  $R^*$  is approximately the same as  $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Similarly, denote  $\mathbf{D} \equiv \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}$  as the set of ReLU neural networks  $D_\phi : \mathbb{R}^d \rightarrow \mathbb{R}$  with parameter  $\phi$ , depth  $\tilde{\mathcal{H}}$ , width  $\tilde{\mathcal{W}}$ , size  $\tilde{\mathcal{S}}$ .

Let  $\{W_i\}_{i=1}^n$  be  $n$  i.i.d random vectors drawn from  $\gamma_d$ . The estimated DDRM is defined by

$$\hat{R}_\theta \in \arg \min_{R_\theta \in \mathbf{R}} \hat{\mathcal{L}}(R_\theta) \quad (9)$$

where  $\hat{\mathcal{L}}(R_\theta) = -\hat{\mathcal{V}}_n[R_\theta(X), Y] + \lambda \hat{\mathbb{D}}_f(\mu_{R_\theta(X)} \| \gamma_d)$ . Here  $\hat{\mathcal{V}}_n[R_\theta(X), Y]$  is an unbiased and consistent estimator of  $\mathcal{V}[R_\theta(X), Y]$  as defined in (5) based on  $\{(R_\theta(X_i), Y_i), i = 1, \dots, n\}$  and

$$\hat{\mathbb{D}}_f(\mu_{R_\theta(X)} \| \gamma_d) = \max_{D_\phi \in \mathbf{D}} \frac{1}{n} \sum_{i=1}^n [D_\phi(R_\theta(X_i)) - f^*(D_\phi(W_i))]. \quad (10)$$

This objective function consists of two terms: (a) the term  $\lambda \hat{\mathcal{V}}_n[R_\theta(X), Y]$  is an unbiased and consistent estimator of  $\lambda \mathcal{V}[R_\theta(X), Y]$ , which is a measure that quantifies the conditional

independence  $X \perp\!\!\!\perp Y | R_{\theta}(X)$ ; (b) the term  $\widehat{\mathbb{D}}_f(\mu_{R_{\theta}(X)} || \gamma_d)$  promotes disentanglement among the components of  $R_{\theta}(X)$  by encouraging  $R_{\theta}(X)$  to be distributed as  $N(0, \mathbf{I}_d)$ . This is the dual form of the  $f$ -GAN loss (Goodfellow et al., 2014; Nowozin et al., 2016). We note that GANs seek to find a map from a reference distribution such as Gaussian to the data space, here we do the reverse and try to find a representation of the data to be distributed like a reference distribution.

## 4 Consistency

We establish the consistency of the estimated DDRM in the sense that the excess risk  $\mathcal{L}(\widehat{R}_{\theta}) - \mathcal{L}(R^*)$  converges to zero, where  $\widehat{R}_{\theta}$  is the deep nonparametric estimator in (9). It is clear that to achieve consistency, it is necessary to require the network parameters to increase as the sample size increases. This is similar to requiring the bandwidth of a nonparametric kernel density estimator to depend on the sample size. There is an extensive literature on how to select the bandwidth parameter in nonparametric density estimation problems. How to choose the structure parameters of a neural network is a more complicated problem. To the best of our knowledge, it has not been systematically studied in the literature. We provide a particular specification below that ensures the consistency of the estimated representation. However, this specification is not necessarily optimal, it only represents our first attempt to tackle this difficult problem.

We make the following basic assumptions about the target parameter and the model.

- (A1) The target representation  $R^*$  is Lipschitz continuous with Lipschitz constant  $L_1$ .
- (A2) For every  $R \in \mathbf{R} \equiv \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ , we assume the density ratio  $r(\mathbf{z}) = \frac{d\mu_{R(X)}}{d\gamma_d}(\mathbf{z})$  to be Lipschitz continuous with Lipschitz constant  $L_2$ , and  $c_1 \leq r(\mathbf{z}) \leq c_2$  for some constants

$$0 < c_1 \leq c_2 < \infty.$$

(A3)  $\text{supp}(\mu_X)$  is contained in a compact set, say  $[-B_1, B_1]^p$  with a finite  $B_1$  and denote its density function as  $f_X(x)$ .  $Y$  is bounded almost surely, say  $\|Y\| \leq C_1$  a.s..

Let  $B_2 = \max\{|f'(c_1)|, |f'(c_2)|\}$  and  $B_3 = \max_{|s| \leq 2B_2} |f^*(s)|$ . For the KL-divergence, we have  $B_2 = \max\{\log c_1, \log c_2\} + 1$  and  $B_3 = \exp(2B_2)$ . We specify the network parameters of the representer  $R_\theta$  and the discriminator  $D_\phi$  as follows.

(N1) Representer network  $\mathbf{R} \equiv \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  parameters: depth  $\mathcal{H} = \mathcal{O}(\log n)$  width  $\mathcal{W} = \mathcal{O}(n^{\frac{p}{2(2+p)}} / \log n)$ , size  $\mathcal{S} = \mathcal{O}(dn^{\frac{p}{2+p}} / \log^4(npd))$ , and  $\|R\|_{L^\infty} \leq 2\|R^*\|_{L^\infty}, \forall R \in \mathbf{R}$ .

(N2) Discriminator network  $\mathbf{D} \equiv \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}$  parameters: depth  $\tilde{\mathcal{H}} = \mathcal{O}(\log n)$ , width  $\tilde{\mathcal{W}} = \mathcal{O}(n^{\frac{d}{2(2+d)}} / \log n)$ , size  $\tilde{\mathcal{S}} = \mathcal{O}(n^{\frac{d}{2+d}} / \log^4(npd))$ , and  $\|D\|_{L^\infty} \leq 2B_2, \forall D \in \mathbf{D}$ .

We again note that these specifications of the network parameters are not necessarily unique or optimal. Our goal here is to provide theoretical support for the proposed method in the sense that there exist networks with the above specifications leading to the consistency of the estimated representation map.

**Theorem 4.1.** *Set  $\lambda = \mathcal{O}(1)$ . Suppose conditions (A1)-(A3) hold and set the network parameters according to (N1)-(N2). Then  $\mathbb{E}_{\{X_i, Y_i, W_i\}_{i=1}^n} [\mathcal{L}(\hat{R}_\theta) - \mathcal{L}(R^*)] \rightarrow 0$ .*

The proof of this theorem is given in the appendix. Conditions (A1) and (A2) are regularity conditions that are often assumed in nonparametric estimation problems. The result established in Theorem 4.1 shows that the learned DDRM achieves asymptotic sufficiency under the conditions (A1) and (A2) and with the specifications (N1) and (N2) for the network parameters.

There have been intensive efforts devoted to understanding the theoretical properties of deep neural network models in recent years. Several stimulating papers have studied the

statistical convergence properties of nonparametric regression using neural networks (Bauer and Kohler, 2019; Schmidt-Hieber, 2020; Farrell et al., 2021; Jiao et al., 2021). There have also been some recent works on the non-asymptotic error bounds of GANs. For example, Zhang et al. (2018) considered the generalization error of GANs. Liang (2020) studied the rates of convergence for learning distributions implicitly with GAN under several forms of the integral probability metrics. Bai et al. (2019) analyzed the estimation error bound of GANs under the Wasserstein distance for a special class of distributions implemented by a generator. Chen et al. (2020) studied the convergence rates of GAN distribution estimators when both the evaluation class and the target density class are Hölder classes.

In the present problem, the objective function (9) is the combination of a loss that is a  $U$ -process indexed by a class of neural networks and a GAN-type loss indexed by two classes of neural networks. This objective function is more complicated than the least squares loss or the GAN loss analyzed in the aforementioned works. Therefore, the problem here is more difficult. To the best of our knowledge, the consistency property of the excess risk of the minimizer of such an objective function has not been analyzed in the literature.

## 5 Computation

Lemma 3.1 implies that training of  $\phi$  with fixed  $\theta$  is to push forward the distribution of  $R(X)$  to the reference distribution  $\gamma_d = N(0, \mathbf{I}_d)$ . For this purpose, we need to estimate an optimal discriminator  $D_\phi$  approximating the optimal dual function  $D(\mathbf{z}) = f'(r(\mathbf{z}))$ , where  $r(\mathbf{z})$  is the ratio for the density of  $\mu_{R_\theta(\mathbf{x})}$  over the density of  $\gamma_d$ . Note that  $f'$  is a strictly increasing function if  $f$  is strictly convex, which is true for all the commonly used divergence measures. Thus the problem of estimating the discriminator is essentially that of estimating the density ratio. Therefore, in our implementation, we utilize the computationally stable

particle method based on gradient flow in probability measure spaces (Gao et al., 2019, 2020). The key idea of this particle method is to seek a sequence of nonlinear but simpler residual maps,  $\mathbb{T}(\mathbf{z}) = \mathbf{z} + s\mathbf{v}(\mathbf{z})$ , where  $s > 0$  is a small step size, pushing the samples from  $\mu_{R_{\boldsymbol{\theta}}(\mathbf{x})}$  to the target distribution  $\gamma_d$  along a velocity fields  $\mathbf{v}(\mathbf{z}) = -\nabla f'(r(\mathbf{z}))$  that most decreases the  $f$ -divergence  $\mathbb{D}_f(\cdot || \gamma_{d^*})$  at  $\mu_{R_{\boldsymbol{\theta}}(X)}$  (Gao et al., 2019). The residual maps can be estimated via deep density-ratio estimation. Specifically, the estimated residual maps take the form  $\mathbb{T}(\mathbf{z}) = \mathbf{z} + s\widehat{\mathbf{v}}(\mathbf{z})$ ,  $\mathbf{z} \in \mathbb{R}^d$ , where  $\widehat{\mathbf{v}}(\mathbf{z}) = -\nabla f'(\widehat{r}(\mathbf{z}))$ . Here  $\widehat{r}(\mathbf{z})$  is an estimated density ratio of the density of  $R_{\boldsymbol{\theta}}(\mathbf{x})$  at the current value of  $\boldsymbol{\theta}$  over the density of the reference distribution. The estimator  $\widehat{r}(\mathbf{z})$  is constructed as follows. Let  $Z_i = R_{\boldsymbol{\theta}}(X_i)$  and generate  $W_i \sim \gamma_d, i = 1, 2, \dots, n$ . We solve

$$\widehat{D}_{\phi} \in \arg \min_{D_{\phi}} \frac{1}{n} \sum_{i=1}^n \{\log[1 + \exp(D_{\phi}(Z_i))] + \log[1 + \exp(-D_{\phi}(W_i))]\} \quad (11)$$

with stochastic gradient descent (SGD). Then the estimated density ratio  $\widehat{r}(\mathbf{z}) = \exp(-\widehat{D}_{\phi}(\mathbf{z}))$ . Here we note that the population version of the loss function in (11) is minimized at  $-\log(r(\mathbf{z}))$ . Therefore,  $\widehat{D}_{\phi}(\mathbf{z})$  in (11) provides a good estimator of  $-\log(r(\mathbf{z}))$ . See Gao et al. (2020) for a detailed description of this particle approach. Here, we use this approach to transform  $Z_i = R_{\boldsymbol{\theta}}(X_i), i = 1, \dots, n$  into Gaussian samples (we still denote them as  $Z_i$ ) directly. Once this is done, we update  $\boldsymbol{\theta}$  via minimizing the loss

$$\frac{1}{n} \sum_{i=1}^n \|R_{\boldsymbol{\theta}}(X_i) - Z_i\|^2 - \lambda \widehat{\mathcal{V}}_n[R_{\boldsymbol{\theta}}(X), Y].$$

We depict the DDR algorithm in the flowchart in Figure 1 and give a detailed description below. **Pseudo-code for the DDR algorithm**

- Input  $\{X_i, Y_i\}_{i=1}^n$ . Tuning parameters:  $s, \lambda, d$ . Sample  $\{W_i\}_{i=1}^n \sim \gamma_d$ .

- *Outer loop for  $\theta$* 
  - *Inner loop (particle method)*
    - \* Let  $Z_i = R_\theta(X_i), i = 1, \dots, n$ .
    - \* Solve  $\hat{D}_\phi \in \arg \min_{D_\phi} \frac{1}{n} \sum_{i=1}^n \{\log[1 + \exp(D_\phi(Z_i))] + \log[1 + \exp(-D_\phi(W_i))]\}$ .
    - \* Define the residual map  $\mathbb{T}(\mathbf{z}) = \mathbf{z} - s \nabla f'(\hat{r}(\mathbf{z}))$  with  $\hat{r}(\mathbf{z}) = \exp(-\hat{D}_\phi(\mathbf{z}))$ .
    - \* Update the particles  $Z_i = \mathbb{T}(Z_i), i = 1, 2, \dots, n$ .
  - *End inner loop*
  - Update  $\theta$  via minimizing  $-\hat{\mathcal{V}}_n[R_\theta(X), Y] + \lambda \sum_{i=1}^n \|R_\theta(X_i) - Z_i\|^2/n$  using SGD.
- *End outer loop*

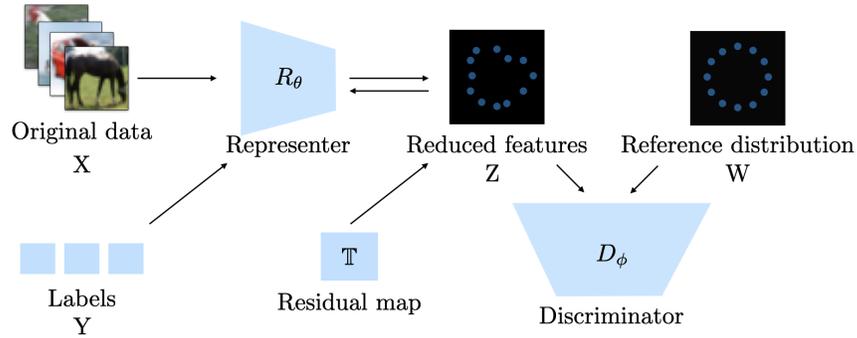


Figure 1: Flow chart for deep dimension reduction (DDR)

## 6 Numerical experiments

We evaluate the performance of DDR using simulated and benchmark real data. Since DDR is not trying to estimate a classifier or a regression function directly, but rather to learn a representation with the desired properties of sufficiency, low-dimensionality and disentanglement, we design the experiments to evaluate the performance of the learned representations based on DDR in terms of prediction when using these representations. The

Table 2: Summary information of DDR and compared methods.  $X^l$  and  $X^u$  represent labeled and unlabeled data, respectively, while  $Y^l$  represents labeled targets.

Method	Name	Input	Supervision	Based Model
DDR	Deep Dimension Reduction	$X^l, Y^l$	Supervised	Neural networks
NN	Neural Networks	$X^l, Y^l$	Supervised	Neural networks
dCorAE	Distance Correlation Autoencoder	$X^l, Y^l$	Supervised	Neural networks
OLS	Ordinary Least Squares	$X^l, Y^l$	Supervised	Linear
SIR	Sliced Inverse Regression	$X^l, Y^l$	Supervised	Linear
SAVE	Sliced Average Variance Estimation	$X^l, Y^l$	Supervised	Linear
GSIR	Generalized Sliced Inverse Regression	$X^l, Y^l$	Supervised	Kernel
GSAVE	Generalized Sliced Average Variance Estimation	$X^l, Y^l$	Supervised	Kernel
Semi-VAE	Semi-supervised Variational Autoencoders	$X^l, Y^l$ and $X^u$	Semi-supervised	Neural networks
InfoVAE	Information Maximizing Variational Autoencoders	$X^l, Y^l$ and $X^u$	Semi-supervised	Neural networks
PCA	Principal Component Analysis	$X^l$ and $X^u$	Unsupervised	Linear
SPCA	Sparse Principal Component Analysis	$X^l$ and $X^u$	Unsupervised	Linear

results demonstrate that a simple classification or regression model using the learned representations performs better than or comparably with the best classification or regression methods using deep neural networks. Details on the network structures and hyperparameters are included in the appendix. Summary information of DDR and compared methods, including the names of methods, their input, learning types, and models of methods, is given in Table 2. Our experiments were conducted on Nvidia DGX Station workstation using a single Tesla V100 GPU unit.

## 6.1 Simulated data

In this subsection, we evaluate DDR on simulated regression and classification problems.

**Regression I.** We generate 10,000 data points from two models:

$$\text{Model (a). } Y = \mathbf{x}_1[0.5 + (\mathbf{x}_2 + 1.5)^2]^{-1} + (1 + \mathbf{x}_2)^2 + \sigma\boldsymbol{\varepsilon}, \text{ where } X \sim N(\mathbf{0}, \mathbf{I}_{20});$$

$$\text{Model (b). } Y = \sin^2(\pi X_1 + 1) + \sigma\boldsymbol{\varepsilon}, \text{ where } X \sim \text{Uniform}[0, 1]^{20}.$$

In both models,  $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{I})$ . We use a 3-layer network with ReLU activation for  $R_\theta$  and a single hidden layer ReLU network for  $D_\phi$ . We compare DDR with four prominent

sufficient dimension reduction methods: sliced inverse regression (SIR) (Li, 1991), sliced average variance estimation (SAVE) (Cook and Weisberg, 1991), generalized sliced inverse regression (GSIR) and generalized sliced average variance estimation (GSAVE) (Lee et al., 2013; Li, 2018). SIR slices the range of  $Y$  and obtains the crude estimation of the inverse regression  $E(X|Y)$ . Then the eigenvectors of the covariance matrix  $Cov(E(X|Y))$  that lie in the central subspace of data can be estimated via weighted PCA. SIR is a first moment method to estimate the central subspace from  $E(X|Y)$ , while SAVE is a second moment method to estimate the space from  $Var(E(X|Y))$  that is primarily used to solve symmetric data problems. Similarly, SAVE also utilizes the weighted PCA to estimate eigenvectors that lie in the central subspace. GSIR and GSAVE are generalized versions of SIR and SAVE, respectively. Both of them estimate central subspace in the reproducing kernel Hilbert space (RKHS) instead of using the covariance matrix in both SIR and SAVE. Also, we compare DDR with two deep learning based methods: neural networks (NN) with least square (LS) loss as the last layer, denoted as NN+LS, and distance correlation autoencoder (dCorAE) (Wang et al., 2018). dCorAE targets at two objectives for both reconstruction and classification, presenting a trade-off between two tasks during training.

Table 3: Average prediction errors and their standard errors (based on 5-fold validation)

Method	Model (a)			Model (b)		
	$\sigma = 0.1$	$\sigma = 0.4$	$\sigma = 0.8$	$\sigma = 0.05$	$\sigma = 0.1$	$\sigma = 0.2$
DDR	<b>0.127 ± .005</b>	<b>0.555 ± .010</b>	<b>1.088 ± .009</b>	<b>0.052 ± .001</b>	<b>0.105 ± .003</b>	<b>0.241 ± .010</b>
NN+LS	0.147 ± .028	0.575 ± .008	1.150 ± .013	0.053 ± .001	0.107 ± .002	0.242 ± .010
dCorAE	0.153 ± .015	0.549 ± .012	1.101 ± .015	0.065 ± .001	0.135 ± .001	0.275 ± .004
SIR	1.484 ± .047	1.599 ± .050	1.712 ± .037	0.252 ± .002	0.268 ± .002	0.323 ± .005
SAVE	1.482 ± .048	1.588 ± .049	1.715 ± .038	0.252 ± .002	0.268 ± .003	0.323 ± .005
GSIR	1.477 ± .047	1.598 ± .050	1.707 ± .039	0.267 ± .004	0.269 ± .004	0.322 ± .006
GSAVE	1.478 ± .048	2.602 ± .079	2.654 ± .041	0.265 ± .003	0.267 ± .004	0.339 ± .006

We fit a linear model with the learned features and the response variable, and report the

prediction errors in Table 3. We see that DDR outperforms SIR, SAVE, GSIR, GSAVE, NN+LS and dCorAE in terms of prediction error.

**Regression II.** We generate 5000 data points from three simulated models:

$$\text{Model (a). } Y = (\mathbf{x}_1 + \mathbf{x}_2)^2 + (1 + \exp(\mathbf{x}_1))^2 + \varepsilon;$$

$$\text{Model (b). } Y = \sin(\pi(\mathbf{x}_1 + \mathbf{x}_2)/10) + \mathbf{x}_1^2 + \varepsilon;$$

$$\text{Model (c). } Y = (\mathbf{x}_1^2 + \mathbf{x}_2^2)^{1/2} \log(\mathbf{x}_1^2 + \mathbf{x}_2^2)^{1/2} + \varepsilon,$$

where  $\varepsilon \perp X$  and  $\varepsilon \sim N(\mathbf{0}, 0.25 \cdot \mathbf{I}_{10})$ . For the distribution of the 10-dimensional predictor  $X$ , we consider three scenarios: Scenario (i):  $X \sim N(\mathbf{0}, \mathbf{I}_{10})$ ; independent Gaussian predictors; Scenario (ii):  $X \sim \frac{1}{3}N(-2 \cdot \mathbf{1}_{10}, \mathbf{I}_{10}) + \frac{1}{3}\text{Uniform}[-1, 1]^{10} + \frac{1}{3}N(2 \cdot \mathbf{1}_{10}, \mathbf{I}_{10})$ , independent non-Gaussian predictors; Scenario (iii):  $X \sim N(\mathbf{0}, 0.3 \cdot \mathbf{I}_{10} + 0.7 \cdot \mathbf{1}_{10}\mathbf{1}_{10}^\top)$ . correlated Gaussian predictors. These models and the distributional scenarios are modified from (Lee et al., 2013; Li, 2018).

Table 4: Average prediction errors (APE), distance correlation (DC), conditional Hilbert-Schmidt independence criterion (HSIC) and their standard errors (based on 5-fold validation)

	Method	Model (a)			Model (b)			Model (c)		
		APE	DC	HSIC	APE	DC	HSIC	APE	DC	HSIC
Scenario (i)	DDR	<b>6.1 ± 3.5</b>	<b>1.0 ± .0</b>	<b>34.7 ± 3.5</b>	<b>0.3 ± .0</b>	<b>1.0 ± .0</b>	<b>49.1 ± 7.7</b>	<b>0.3 ± .0</b>	<b>0.9 ± .0</b>	<b>36.4 ± 2.8</b>
	GSIR	28.3 ± 7.5	0.2 ± .0	64.4 ± 2.9	1.4 ± .1	0.1 ± .0	134.8 ± 13.5	0.8 ± .0	0.2 ± .0	67.4 ± 3.9
	GSAVE	28.3 ± 7.4	0.1 ± .0	72.1 ± 4.1	1.4 ± .0	0.1 ± .0	175.9 ± 6.4	0.8 ± .0	0.2 ± .0	66.5 ± 2.8
Scenario (ii)	DDR	<b>183.1 ± 98.7</b>	<b>0.9 ± .1</b>	45.1 ± 3.3	<b>0.4 ± .1</b>	<b>1.0 ± .0</b>	<b>16.8 ± 2.0</b>	<b>0.3 ± .1</b>	<b>1.0 ± .0</b>	<b>8.6 ± 0.7</b>
	GSIR	664.6 ± 38.1	0.1 ± .0	<b>43.9 ± 2.6</b>	3.3 ± .2	0.1 ± .0	27.6 ± 1.7	1.5 ± .1	0.6 ± .0	14.6 ± 0.8
	GSAVE	662.0 ± 38.0	0.0 ± .0	48.0 ± 2.7	3.2 ± .2	0.2 ± .0	32.0 ± 1.5	2.4 ± .0	0.0 ± .0	15.5 ± 0.6
Scenario (iii)	DDR	<b>12.5 ± 11.1</b>	<b>0.8 ± .3</b>	<b>37.0 ± 4.7</b>	<b>0.3 ± .0</b>	<b>1.0 ± .0</b>	<b>48.6 ± 7.1</b>	<b>0.3 ± .1</b>	<b>0.9 ± .1</b>	36.8 ± 5.9
	GSIR	32.2 ± 6.1	0.2 ± .1	61.4 ± 6.0	1.0 ± .1	0.2 ± .0	51.3 ± 51.3	0.6 ± .0	0.6 ± .0	<b>25.8 ± 25.6</b>
	GSAVE	31.8 ± 6.4	0.2 ± .1	60.0 ± 4.4	1.0 ± .0	0.3 ± .0	119.3 ± 8.8	0.6 ± .0	0.7 ± .0	54.3 ± 3.8

We compare DDR with generalized sliced inverse regression (GSIR) and generalized sliced average variance estimation (GSAVE) (Lee et al., 2013; Li, 2018). In DDR, we adopt a 4-layer network for  $R_\theta$  and a 3-layer network for  $D_\phi$  with Leaky ReLU activation. For

all methods, we fit a linear model with the learned features and the response variable, and report the prediction error, distance correlation between representation and the response variable, and conditional Hilbert-Schmidt independence criterion (HSIC) (Fukumizu et al., 2008). The results are presented in Table 4. The representations learned with DDR present higher distance correlations with the response and lower conditional HSICs, suggesting that DDR is capable of better capturing data information and conditional independence property than other methods. Moreover, DDR significantly outperforms GSIR and GSAVE in terms of prediction errors in all scenarios.

**Classification.** We visualize the learned features of DDR on three simulated datasets. We first generate (1) 2-dimensional concentric circles from two classes as in Figure 2 (a); (2) 2-dimensional moons data from two classes as in Figure 2 (e); (3) 3-dimensional Gaussian data from six classes as in Figure 2 (i). In each dataset, we generate 5,000 data points for each class. We next map the data into 100-dimensional space using matrices with entries i.i.d  $\text{Unifrom}([0, 1])$ . Finally, we apply DDR to these 100-dimensional datasets to learn 2-dimensional features. We use a 10-layer dense convolutional network (DenseNet) (Huang et al., 2017a) as  $R_{\theta}$  and a 4-layer network as  $D_{\phi}$ . We display the evolutions of the learned 2-dimensional features by DDR in Figure 2. For ease of visualization, we push all the distributions onto the uniform distribution on the unit circle, which is done by normalizing the standard Gaussian random vectors to length one. Clearly, the learned features for different classes in the examples are well disentangled.

## 6.2 Real datasets

We benchmark DDR on a variety of real datasets from both regression and classification problems. Summary information of those datasets used in the analysis is given in Table 5.

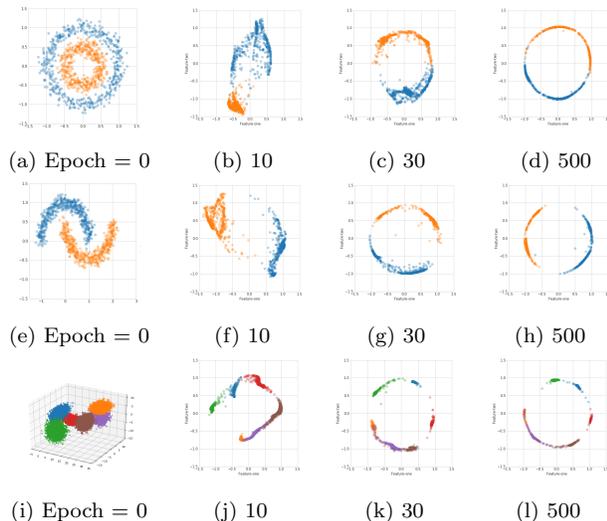


Figure 2: Evolving learned features at Epoch = 0, 10, 30, and 500. The first, second and third rows show concentric circles, moons and 3D Gaussian datasets, respectively.

Table 5: Summary information for real datasets.

Dataset	Feature size	Training size	Test size	Task
YearPredictionMSD	90	412,276	103,069	Regression
Pole-Telecommunication	48	12,000	3,000	Regression
MNIST	$28 \times 28 \times 1$	$60k$	$10k$	Classification with 10 categories
Kuzushiji-MNIST	$28 \times 28 \times 1$	$60k$	$10k$	Classification with 10 categories
FashionMNIST	$28 \times 28 \times 1$	$60k$	$10k$	Classification with 10 categories
CIFAR-10	$32 \times 32 \times 3$	$50k$	$10k$	Classification with 10 categories
CIFAR-100	$32 \times 32 \times 3$	$50k$	$10k$	Classification with 100 categories

**Regression.** We benchmark the prediction performance of regression models using the representations learned based on DDR. Here, we use the YearPredictionMSD dataset<sup>1</sup> and the Pole-Telecommunication dataset<sup>2</sup>. The YearPredictionMSD dataset contains 515,345 observations with 90 predictors. The problem is to predict the year of song release. The

<sup>1</sup>The YearPredictionMSD dataset is available at <https://archive.ics.uci.edu/ml/datasets/YearPredictionMSD>.

<sup>2</sup>The Pole-Telecommunication dataset is available at <https://www.dcc.fc.up.pt/~ltorgo/Regression/DataSets.html>.

Pole-Telecommunication dataset consists of 15,000 observations with 48 predictors for determining the placement of antennas. We randomly split the data into five folds to evaluate the prediction performance using 5-fold cross validation. We employ a 3-layer network for both  $D_\phi$  and  $R_\theta$  on the YearPredictionMSD dataset; a 2-layer network for  $D_\phi$  and a 4-layer network  $R_\theta$  are adopted on the Pole-Telecommunication dataset. In comparison, we conduct a nonlinear regression using neural networks (NN) with a least squares (LS) loss in the last layer, denoted as NN + LS. That is, we do not impose any desired characteristics for the learned representations in the pultimate layer. Note that, for both DDR and NN+LS, we use the same networks to learn representative features. We also consider the popular dimension reduction methods, including principal component analysis (PCA) and sparse principal component analysis (SPCA), to obtain data representation. For the comparison with supervised dimension reduction methods, we consider SIR and SAVE and the deep learning based sufficient dimension reduction method dCorAE. For those methods, we first obtain the estimated representative features and fit a linear regression model of the response on the learned representations. The average prediction errors and their standard errors based on DDR, NN+LS, dCorAE, PCA, SPCA, SIR, SAVE and the ordinary least squares (OLS) regression with the original data are reported in Tables 6 and 7. DDR outperforms other methods in terms of prediction accuracy.

Table 6: Prediction error  $\pm$  standard error: YearPredictionMSD dataset

Methods	$d = 10$	$d = 20$	$d = 30$	$d = 40$
DDR	<b>8.8 <math>\pm</math> 0.1</b>	<b>8.9 <math>\pm</math> 0.1</b>	<b>8.9 <math>\pm</math> 0.1</b>	<b>8.8 <math>\pm</math> 0.1</b>
dCorAE	8.9 $\pm$ 0.1	9.0 $\pm$ 0.1	9.2 $\pm$ 0.1	8.9 $\pm$ 0.1
NN+LS	9.2 $\pm$ 0.1	9.3 $\pm$ 0.1	9.2 $\pm$ 0.1	9.2 $\pm$ 0.1
SPCA	10.6 $\pm$ 0.1	10.4 $\pm$ 0.1	9.6 $\pm$ 0.1	10.2 $\pm$ 0.1
PCA	10.6 $\pm$ 0.1	10.4 $\pm$ 0.1	10.3 $\pm$ 0.1	10.2 $\pm$ 0.1
SIR	9.6 $\pm$ 0.1	9.6 $\pm$ 0.1	9.6 $\pm$ 0.1	9.6 $\pm$ 0.1
SAVE	10.3 $\pm$ 0.1	9.7 $\pm$ 0.1	9.6 $\pm$ 0.1	9.6 $\pm$ 0.1
OLS		9.6 $\pm$ 0.1		

Table 7: Prediction error  $\pm$  standard error: Pole-Telecommunication dataset

Methods	$d = 5$	$d = 10$	$d = 15$	$d = 20$
DDR	<b>2.1 <math>\pm</math> 0.2</b>	<b>2.1 <math>\pm</math> 0.1</b>	<b>2.2 <math>\pm</math> 0.1</b>	<b>2.2 <math>\pm</math> 0.2</b>
dCorAE	3.1 $\pm$ 0.1	3.1 $\pm$ 0.3	3.1 $\pm$ 0.2	3.0 $\pm$ 0.1
NN +LS	3.0 $\pm$ 0.5	3.1 $\pm$ 1.1	2.7 $\pm$ 0.8	3.2 $\pm$ 0.6
SPCA	40.3 $\pm$ 0.3	40.1 $\pm$ 0.3	30.5 $\pm$ 0.2	30.5 $\pm$ 0.2
PCA	40.3 $\pm$ 0.3	40.1 $\pm$ 0.3	30.5 $\pm$ 0.2	30.5 $\pm$ 0.2
SIR	30.4 $\pm$ 0.1	30.4 $\pm$ 0.1	30.5 $\pm$ 0.1	30.5 $\pm$ 0.1
SAVE	31.2 $\pm$ 0.3	30.5 $\pm$ 0.1	30.5 $\pm$ 0.1	30.5 $\pm$ 0.1
OLS		30.5 $\pm$ 0.2		

**Classification I.** We benchmark the classification performance of DDR using MNIST (LeCun et al., 2010), FashionMNIST (Xiao et al., 2017), CIFAR-10, and CIFAR-100 (Krizhevsky and Hinton, 2009) datasets against some existing methods, including neural networks (NN) with cross entropy (CN) loss as the last layer, denoted as CNN, and distance correlation autoencoder (dCorAE) (Wang et al., 2018). With CNN, we use the feature extractor by dropping the last layer for the CN loss of the NN trained for classification as networks. Note that, for both DDR and CNN, we apply the same networks to learn representations.

The MNIST and FashionMNIST datasets consist of  $60k$  and  $10k$  grayscale images with  $28 \times 28$  pixels for training and testing, respectively, while the CIFAR-10 and CIFAR-100 datasets contain  $50k$  and  $10k$  colored images with  $32 \times 32$  pixels for training and testing, respectively. The representer network  $R_\theta$  contains 20 layers for MNIST data and 100 layers for CIFAR-10 data.

To fully utilize computational resources and improve classification accuracy, we further combine DDR with the CN loss, denoted as DDR+CN, by applying the transfer learning technique (Torrey and Shavlik, 2010; Pan and Yang, 2009; Tan et al., 2018) on CIFAR-10 and CIFAR-100. Data structures for both CIFAR-10 and ImageNet are the same (with three channels), which makes the use of transfer learning straightforward by leveraging the

pretrained model of ImageNet. The pretrained WideResnet-101 model (Zagoruyko and Komodakis, 2016) on the ImageNet dataset with Spinal FC (Kabir et al., 2020) is chosen for  $R_\theta$ . In our experiments for transfer learning, we first train the WideResnet model on ImageNet. We then use the parameters of the pretrained neural network as the initialization parameters to train CIFAR-10. In contrast to transfer learning, the initialization parameters of learning from scratch are random. The discriminator network  $D_\phi$  is a 4-layer network. The architecture of  $R_\theta$  and most hyperparameters are shared across all four methods - DDR, CNN, DDR+CN and dCorAE (Wang et al., 2018). Finally, we use the  $k$ -nearest neighbor ( $k = 5$ ) classifier on the learned features for all methods.

As shown in Table 8, the classification accuracies of DDR for MNIST and FashionMNIST are better than or comparable with those of CNN and dCorAE. As shown in Table 9, the classification accuracy of DDR using the CN loss outperforms that of CNN on CIFAR-10 and CIFAR-100. We also calculate the estimated distance correlation (DC) between the learned features and their labels. Figure 3 shows the values of DC for MNIST, FashionMNIST and CIFAR-10 data. Higher DC values mean that the learned features are of higher quality. DDR and DDR+CN achieves higher DC values.

Because both GSIR and GSAVE require computation of  $n \times n$  kernel matrices, which is computationally prohibitive when  $n = 10,000$  to  $60,000$  and  $p \approx 1,000$ , it does not allow us to apply both methods to analyze datasets from MNIST, FashionMNIST, CIFAR-10 and CIFAR-100.

**Classification II.** To compare the performance of DDR with semi-supervised methods, we benchmark DDR on MNIST dataset with varying amounts of labeled data for training. In detail, we consider some widely used semi-supervised learning methods, including semi-supervised variational autoencoders (Semi-VAE) (Kingma et al., 2014) and information maximizing variational autoencoders (InfoVAE) with the semi-supervised setting (Kingma

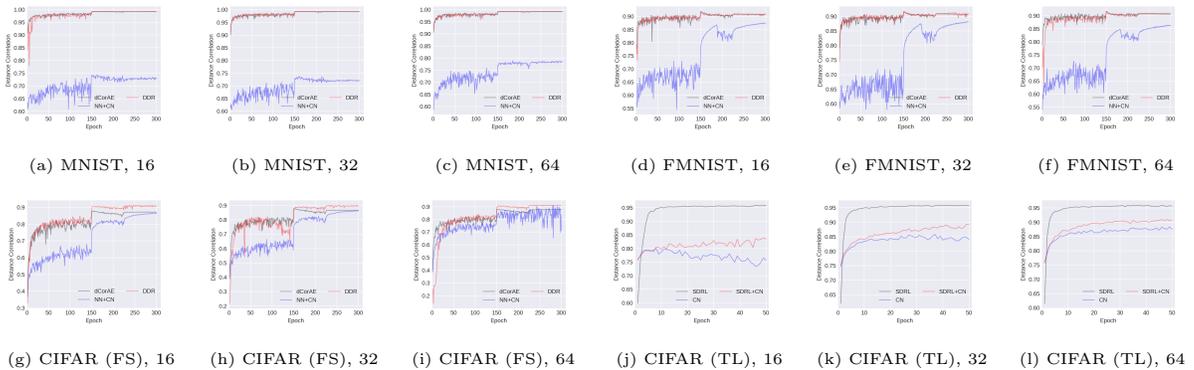


Figure 3: The distance correlations of labels with learned features based on DDR, CNN and dCorAE with  $d = 16, 32$  and  $64$  for MNIST, FashionMNIST (FMNIST) and CIFAR-10 (CIFAR) data (FS: from scratch; TL: transfer learning).

Table 8: Classification accuracy for MNIST and FashionMNIST. In the table, dCor=dCorAE.

$d$	MNIST			FashionMNIST		
	DDR	dCor	CNN	DDR	dCor	CNN
$d = 16$	99.41	99.58	99.39	<b>94.44</b>	94.18	94.21
$d = 32$	<b>99.61</b>	99.54	99.45	94.18	93.89	94.41
$d = 64$	<b>99.56</b>	99.53	99.49	94.13	94.24	94.38

Table 9: Classification accuracy for CIFAR-10 and CIFAR-100 data.

$d$	CIFAR-10						CIFAR-100			
	Learning from scratch			Transfer learning			Transfer learning			
	dCorAE	CNN	DDR	CNN	DDR	DDR+CN	$d$	dCorAE	CNN	DDR+CN
$d = 16$	94.15	94.21	<b>94.29</b>	97.44	97.52	<b>97.68</b>	$d = 200$	85.39	86.29	<b>86.36</b>
$d = 32$	94.18	94.92	94.58	97.79	97.33	<b>97.96</b>	$d = 300$	85.57	85.95	<b>86.04</b>
$d = 64$	94.66	95.09	94.46	97.90	97.49	<b>97.91</b>	$d = 400$	85.55	86.21	<b>86.30</b>

et al., 2014; Zhao et al., 2019), and the supervised method, CNN. InfoVAE utilizes all training images to learn representations, and then trains the  $k$ -nearest neighbor ( $k = 5$ ) classifier with the learned representations and partially known labels. All four methods

share the same network architecture for 50-dimensional learning representation. We adopt the double-hidden-layer MLP networks, with 600 neurons for each layer, the softplus activation function, and the Adam optimizer. For both Semi-VAE and InfoVAE, we apply a semi-supervised setting to analyze all  $60k$  images with the varying number of labeled images as the training data and validate the performance using  $10k$  test data. But for both DDR and CNN, we apply a supervised setting to analyze only data with labels and discard the rest of images in the training set.

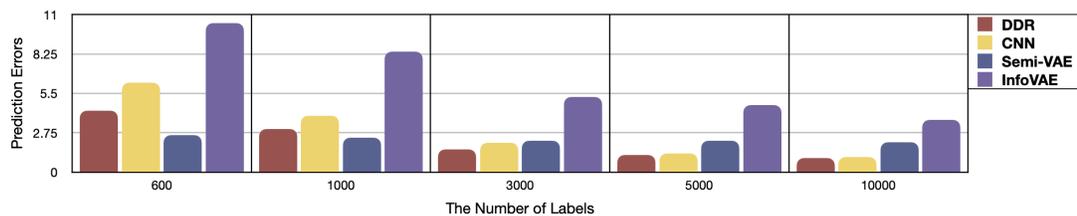


Figure 4: The classification errors comparison on varying amounts of labeled MNIST data.

The prediction accuracy using images with the varying number of labels, from 600 to 10,000, is shown in Fig. 4. For supervised methods, we observe that the accuracy of DDR outperforms that of CNN for the varying number of images with labels. When the proportion of images with labels is low, Semi-VAE performs the best by fully utilizing all  $60k$  training images. However, the accuracy of Semi-VAE does not improve over the increasing proportion of images with labels. This is because the objective function of Semi-VAE is primary to maximize the lower bound of the joint likelihood rather than the classification loss. In all, compared with semi-supervised learning, DDR uses a small amount of data with labels in training, but achieves better classification accuracy when the number of images with labels are larger than 1,000.

## 7 Conclusion and future work

Since the framework of supervised dimension reduction was first introduced over twenty years ago by Li (1991), there have been a variety of important methods developed in this framework. However, most of the existing works focused on finding linear representations of the input data. There are two important reasons why it has been difficult to develop nonparametric supervised representation learning approaches. First, supervised representation learning is more difficult than and fundamentally different from supervised learning. Indeed, it is challenging to formulate a clear and simple objective function for supervised representation learning in the first place. This is in clear contrast to supervised learning, whose objective is clear-cut. For example, in classification, the objective is to minimize the misclassification rate or surrogate objective function; in regression, a least squares criterion for the fitting error is usually used. However, how to construct an objective function for nonparametric supervised representation learning in a principled way has remained an open question (Bengio et al., 2013; Alain and Bengio, 2016). Second, it is difficult to apply standard techniques for nonparametric estimation such as smoothing, splines and kernel methods in multi-dimensional problems. They are either not flexible enough for providing accurate adaptive and data-driven based approximation to multi-dimensional functions or are computational prohibitive with high-dimensional data.

In this work, we propose a nonparametric DDR approach to achieving a good data representation for supervised learning with certain desired characteristics including sufficiency, low-dimensionality and disentanglement. We estimate the representation map nonparametrically by taking advantage of the powerful capabilities of deep neural networks in approximating multi-dimensional functions. The proposed DDR is validated via comprehensive numerical experiments and real data analysis in the context of regression

and classification.

Several questions deserve further study. First, it would be interesting to consider other measures of conditional independence such as conditional covariance operators on reproducing kernel Hilbert spaces (Fukumizu et al., 2009) and heteroscedastic conditional variance operator on Hilbert spaces (Lee et al., 2013). It is also possible to use mutual information for measuring conditional independence (Suzuki and Sugiyama, 2013), although with this measure the loss function itself needs to be estimated. It would be interesting to develop algorithms and theoretical understanding of these criteria and evaluate the relative performance of the learned representations based on these different conditional independence measures.

We used the standard Gaussian as the reference distribution for DDRM to promote disentanglement of the representation. Another convenient choice is the uniform distribution on the unit cube. It is worth examining whether there is any difference in the performance of DDR with different reference distributions. Another question is how to determine the dimension of the learned representation. This is also an important problem in linear SDR. Since the purpose of dimension reduction is often to build prediction models in high-dimensional settings, the problem of determining the dimension of the representation can be best addressed based on cross validation or related data-driven methods in the model building phase.

Last but not least, due to the non-uniqueness of the target, it is challenging to provide the consistency of the estimated nonlinear dimension reduction map. It will be interesting to explore this property in the future work.

## 8 Appendix

In the appendix, we show additional experiments using DDR and InfoVAE. In addition, we provide the implementation details about numerical settings, network structures, SGD optimizers, and hyper-parameters used in the numerical studies. We also give the detailed proofs of Lemmas 2.1-3.2 and Theorem 4.1.

### 8.1 Additional experiments

We make comparisons of DDR and InfoVAE on MNIST, Kuzushiji-MNIST and Fashion-MNIST. Because InfoVAE is an unsupervised method based on autoencoders with maximum mean discrepancy (MMD) loss for constrained representation, we use the semi-supervised setting of InfoVAE (Kingma et al., 2014; Zhao et al., 2019) in our experiments. We utilize the double-hidden-layer MLP networks with the softplus activation function, and the Adam optimizer.

Table 10: MMD metric and classification error for MNIST, Kuzushiji-MNIST and Fashion-MNIST.

		MMD (%)				Classification error (%)				
		#label	600	1000	3000	5000	600	1000	3000	5000
MNIST	DDR		<b>0.0249</b>	<b>0.0242</b>	0.0367	0.127	<b>9.92</b>	<b>6.70</b>	<b>2.80</b>	<b>2.03</b>
	InfoVAE		0.0272	0.0273	0.0273	0.0273	13.86	11.27	7.42	6.29
Kuzushiji-MNIST	DDR		0.0214	0.0242	0.0287	0.0247	<b>27.24</b>	<b>22.75</b>	<b>14.04</b>	<b>9.63</b>
	InfoVAE		0.0202	0.0202	0.0202	0.0202	34.58	28.57	18.89	15.19
Fashion-MNIST	DDR		<b>0.0278</b>	<b>0.0344</b>	0.0475	0.0474	<b>21.25</b>	<b>18.36</b>	<b>16.40</b>	<b>15.15</b>
	InfoVAE		0.0404	0.0404	0.0404	0.0404	21.47	20.27	17.66	16.67

We compare DDR with InfoVAE in terms of classification error and MMD metric, as shown in Table 10. We observe that DDR and InfoVAE are comparable to learn a standard Gaussian distribution, but DDR outperforms InfoVAE across all settings in terms of the

classification accuracy. In DDR, there exists a trade-off between two objectives: learn a Gaussian distribution and improve the performance of downstream tasks. From the experimental results, we conclude that DDR can achieve satisfactory prediction accuracy with mild constraints on the representation distributions.

## 8.2 Experimental details

### 8.2.1 Simulation studies

The values of the hyper-parameters for the simulated experiments are given in Table 11, where  $\lambda$  is the penalty parameter,  $d$  is the dimension of the SDRM,  $n$  is the mini-batch size in SGD,  $T_1$  is the number of inner loops to push forward particles  $\mathbf{z}_i$ ,  $T_2$  is the number of outer loops for training  $R_\theta$ , and  $s$  is the step size to update particles. For the regression models, the neural network architectures are shown in Table 12.

As shown in Table 13, a multilayer perceptron (MLP) is utilized for the neural structure  $D_\phi$  in the classification problem. The detailed architecture of 10-layer dense convolutional network (DenseNet) (Huang et al., 2017b; Amos and Kolter) deployed for  $R_\theta$  is shown in Table 14. For all the settings, we adopted the Adam (Kingma and Ba, 2014) optimizer with an initial learning rate of 0.001 and weight decay of 0.0001.

Table 11: Hyper-parameters for simulated examples, where  $s$  varies according to epoch

Task	$\lambda$	$d$	$n$	$T_1$	$T_2$	$s$		
						0-150	151-225	226-500
Regression	1.0	2 or 1	64	1	500	3.0	2.0	1.0
Classification	1.0	2	64	1	500	2.0	1.5	1.0

Table 12: MLP architectures for  $D_\phi$  and  $R_\theta$  in regression

Layers	$D_\phi$		$R_\theta$	
	Details	Output size	Details	Output size
Layer 1	Linear, LeakyReLU	16	Linear, LeakyReLU	16
Layer 2	Linear	1	Linear, LeakyReLU	8
Layer 3			Linear	$d$

Table 13: MLP architecture for  $D_\phi$  in the simulated classification examples and the benchmark classification datasets

Layers	Details	Output size
Layer 1	Linear, LeakyReLU	64
Layer 2	Linear, LeakyReLU	128
Layer 3	Linear, LeakyReLU	64
Layer 4	Linear	1

Table 14: DenseNet architecture for  $R_\theta$  in the simulated classification examples

Layers	Details	Output size
Convolution	$3 \times 3$ Conv	$24 \times 20 \times 20$
Dense Block 1	$\begin{bmatrix} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{bmatrix} \times 1$	$36 \times 20 \times 20$
Transition Layer 1	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$30 \times 10 \times 10$
Dense Block 2	$\begin{bmatrix} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{bmatrix} \times 1$	$18 \times 10 \times 10$
Transition Layer 2	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$15 \times 5 \times 5$
Dense Block 3	$\begin{bmatrix} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{bmatrix} \times 1$	$27 \times 5 \times 5$
Pooling	BN, ReLU, $5 \times 5$ Average Pool, Reshape	27
Fully connected	Linear	2

### 8.2.2 Real datasets

**Regression:** In the regression problems, hyper-parameters are presented in Table 15. The Adam optimizer with an initial learning rate of 0.001 and weight decay of 0.0001 is adopted. The MLP architectures of  $D_\phi$  and  $R_\theta$  for the YearPredictionMSD data are shown in Table 16 and for the Pole-Telecommunication data are shown in Table 17.

Table 15: Hyper-parameters for real regression datasets

Dataset	$\lambda$	$d$	$n$	$T_1$	$T_2$	$s$
YearPredictionMSD	1.0	10, 20, 30, 40	64	1	500	1.0
Pole-Telecommunication	1.0	5, 10, 15, 20	64	1	200	1.0

Table 16: MLP architectures for  $D_\phi$  and  $R_\theta$  for YearPredictionMSD data

Layers	$D_\phi$		$R_\theta$	
	Details	Output size	Details	Output size
Layer 1	Linear, LeakyReLU	32	Linear, LeakyReLU	32
Layer 2	Linear, LeakyReLU	8	Linear, LeakyReLU	8
Layer 3	Linear	1	Linear	$d$

Table 17: MLP architectures for  $D_\phi$  and  $R_\theta$  for Pole-Telecommunication data

Layers	$D_\phi$		$R_\theta$	
	Details	Output size	Details	Output size
Layer 1	Linear, LeakyReLU	8	Linear, LeakyReLU	16
Layer 2	Linear	$d$	Linear, LeakyReLU	32
Layer 3			Linear, LeakyReLU	8
Layer 4			Linear	$d$

**Classification:** We again use Adam as the SGD optimizers for both  $D_\phi$  and  $R_\theta$ . Specifically, learning rate of 0.001 and weight decay of 0.0001 are used for  $D_\phi$  in all datasets and for  $R_\theta$  on MNIST (LeCun et al., 2010). We customized the SGD optimizers with momentum at 0.9, weight decay at 0.0001, and learning rate  $\rho$  in Table 19. For the transfer learning of CIFAR-10, we use customized SGD optimizer with initial learning rate of 0.001 and momentum of 0.9 for  $R_\theta$ . For FashionMNIST (Xiao et al., 2017) and CIFAR-10 (Krizhevsky et al., 2012), MLP architectures of the discriminator network  $D_\phi$  for MNIST, FashionMNIST and CIFAR-10 are given in Table 13. The 20-layer DenseNet networks shown in Table 20 were utilized for  $R_\theta$  on the MNIST dataset, while the 100-layer DenseNet networks shown in Table 21 and 22 are fitted for  $R_\theta$  on FashionMNIST and CIFAR-10.

Table 18: Hyper-parameters for the classification benchmark datasets

Dataset	$\lambda$	$d$	$n$	$T_1$	$T_2$	$s$
MNIST	1.0	16, 32, 64	64	1	300	0.1
FashionMNIST	1.0	16, 32, 64	64	1	300	1.0
CIFAR-10	1.0	16, 32, 64	64	1	300	1.0
CIFAR-10 (transfer learning)	0.01	16, 32, 64	64	1	50	1.0

Table 19: Learning rate  $\rho$  varies during training.

Epoch	0-150	151-225	226-300
$\rho$	0.1	0.01	0.001

Table 20: Architecture for MNIST, reduced feature size is  $d$

Layers	Details	Output size
Convolution	$3 \times 3$ Conv	$24 \times 28 \times 28$
Dense Block 1	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 2$	$48 \times 28 \times 28$
Transition Layer 1	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$24 \times 14 \times 14$
Dense Block 2	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 2$	$48 \times 14 \times 14$
Transition Layer 2	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$24 \times 7 \times 7$
Dense Block 3	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 2$	$48 \times 7 \times 7$
Pooling	BN, ReLU, $7 \times 7$ Average Pool, Reshape	48
Fully connected	Linear	$d$

Table 21: Architecture for FashionMNIST, reduced feature size is  $d$

Layers	Details	Output size
Convolution	$3 \times 3$ Conv	$24 \times 28 \times 28$
Dense Block 1	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 16$	$216 \times 28 \times 28$
Transition Layer 1	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$108 \times 14 \times 14$
Dense Block 2	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 16$	$300 \times 14 \times 14$
Transition Layer 2	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$150 \times 7 \times 7$
Dense Block 3	$\begin{bmatrix} \text{BN}, 1 \times 1 \text{ Conv} \\ \text{BN}, 3 \times 3 \text{ Conv} \end{bmatrix} \times 16$	$342 \times 7 \times 7$
Pooling	BN, ReLU, $7 \times 7$ Average Pool, Reshape	342
Fully connected	Linear	$d$

Table 22: Architecture for CIFAR-10, reduced feature size is  $d$

Layers	Details	Output size
Convolution	$3 \times 3$ Conv	$24 \times 32 \times 32$
Dense Block 1	$\left[ \begin{array}{l} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{array} \right] \times 16$	$216 \times 32 \times 32$
Transition Layer 1	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$108 \times 16 \times 16$
Dense Block 2	$\left[ \begin{array}{l} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{array} \right] \times 16$	$300 \times 16 \times 16$
Transition Layer 2	BN, ReLU, $2 \times 2$ Average Pool, $1 \times 1$ Conv	$150 \times 8 \times 8$
Dense Block 3	$\left[ \begin{array}{l} \text{BN, } 1 \times 1 \text{ Conv} \\ \text{BN, } 3 \times 3 \text{ Conv} \end{array} \right] \times 16$	$342 \times 8 \times 8$
Pooling	BN, ReLU, $8 \times 8$ Average Pool, Reshape	342
Fully connected	Linear	$d$

### 8.3 Proofs

In this section, we prove Lemmas 2.1 and 3.1, and Theorems 3.2 and 4.1.

#### 8.3.1 Proof of Lemma 2.1

*Proof.* By assumption  $\mu$  and  $\gamma_d$  are both absolutely continuous with respect to the Lebesgue measure. The desired result holds since it is a special case of the well known results on the existence of optimal transport (Brenier, 1991; McCann, 1995), see also Theorem 1.28 on page 24 of (Philippis, 2013) for details.  $\square$

#### 8.4 Proof of Lemma 3.1

*Proof.* Our proof follows Keziou (2003). Since  $f(t)$  is convex, then  $\forall t \in \mathbb{R}$ , we have  $f(t) = f^{**}(t)$ , where

$$f^{**}(t) = \sup_{s \in \mathbb{R}} \{st - f^*(s)\}$$

is the Fenchel conjugate of  $f^*$ . By Fermat's rule, the maximizer  $s^*$  satisfies

$$t \in \partial f^*(s^*),$$

i.e.,

$$s^* \in \partial f(t)$$

Plugging the above display with  $t = \frac{d\mu_Z}{d\gamma}(x)$  into the definition of  $f$ -divergence, we derive (6). □

## 8.5 Proof of Theorem 3.2

*Proof.* Without loss of generality, we assume  $d = 1$ . For  $R^*$  satisfying (3) and any  $R \in \mathcal{R}$ , we have  $R = \rho_{(R,R^*)}R^* + \varepsilon_R$ , where  $\rho_{(R,R^*)}$  is the correlation coefficient between  $R$  and  $R^*$ ,  $\varepsilon_R = R - \rho_{(R,R^*)}R^*$ . It is easy to see that  $\varepsilon_R \perp R^*$  and thus  $Y \perp \varepsilon_R$ . As  $(\rho_{(R,R^*)}R^*, Y)$  is independent of  $(\varepsilon_R, 0)$ , then by Theorem 3 of Székely and Rizzo (2009)

$$\begin{aligned} \mathcal{V}[R, \mathbf{y}] &= \mathcal{V}[\rho_{(R,R^*)}R^* + \varepsilon_R, \mathbf{y}] \leq \mathcal{V}[\rho_{(R,R^*)}R^*, \mathbf{y}] + \mathcal{V}(\varepsilon_R, 0) \\ &= \mathcal{V}[\rho_{(R,R^*)}R^*, \mathbf{y}] = |\rho_{(R,R^*)}| \mathcal{V}[R^*, \mathbf{y}] \\ &\leq \mathcal{V}[R^*, \mathbf{y}]. \end{aligned}$$

As  $R(\mathbf{x}) \sim \mathcal{N}(0, 1)$  and  $R^*(\mathbf{x}) \sim \mathcal{N}(0, 1)$ , then  $\mathbb{D}_f(\mu_{R(\mathbf{x})} \| \gamma_d) = \mathbb{D}_f(\mu_{R^*(\mathbf{x})} \| \gamma_d) = 0$ , and

$$\mathcal{L}(R) - \mathcal{L}(R^*) = \mathcal{V}[R^*, \mathbf{y}] - \mathcal{V}[R, \mathbf{y}] \geq 0.$$

The proof is completed. □

## 8.6 Proof of Theorem 4.1

Recall that  $B_2 = \max\{|f'(c_1)|, |f'(c_2)|\}$ ,  $B_3 = \max_{|s| \leq 2B_2} |f^*(s)|$ . We set the network parameters of the representer  $R_\theta$  and the discriminator  $D_\phi$  as follows.

(N1) Representer network  $\mathbf{R} \equiv \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  parameters: depth  $\mathcal{H} = \mathcal{O}(\log n)$  width  $\mathcal{W} = \mathcal{O}(n^{\frac{p}{2(2+p)}} / \log n)$ , size  $\mathcal{S} = \mathcal{O}(dn^{\frac{p}{2+p}} / \log^4(npd))$ , and  $\|R\|_{L^\infty} \leq \mathcal{B} = 2\|R^*\|_{L^\infty}, \forall R \in \mathbf{R}$ .

(N2) Discriminator network  $\mathbf{D} \equiv \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}$  parameters: depth  $\tilde{\mathcal{H}} = \mathcal{O}(\log n)$ , width  $\tilde{\mathcal{W}} = \mathcal{O}(n^{\frac{d}{2(2+d)}} / \log n)$ , size  $\tilde{\mathcal{S}} = \mathcal{O}(n^{\frac{d}{2+d}} / \log^4(npd))$ , and  $\|D\|_{L^\infty} \leq 2B_2, \forall D \in \mathbf{D}$ .

Before getting into the details of the proof of Theorem 4.1, we first give an outline of the basic structure of the proof.

Without loss of generality, we assume that  $\lambda = 1$  and  $m = 1$ , i.e.  $\mathbf{y} \in \mathbb{R}$ . For any  $\bar{R} \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ , we have,

$$\begin{aligned} \mathcal{L}(\hat{R}_\theta) - \mathcal{L}(R^*) &= \mathcal{L}(\hat{R}_\theta) - \hat{\mathcal{L}}(\hat{R}_\theta) + \hat{\mathcal{L}}(\hat{R}_\theta) - \hat{\mathcal{L}}(\bar{R}) + \hat{\mathcal{L}}(\bar{R}) - \mathcal{L}(\bar{R}) + \mathcal{L}(\bar{R}) - \mathcal{L}(R^*) \\ &\leq 2 \sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(R) - \hat{\mathcal{L}}(R)| + \inf_{\bar{R} \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(\bar{R}) - \mathcal{L}(R^*)|, \end{aligned} \quad (12)$$

where we use the definition of  $\hat{R}_\theta$  in (9) and the feasibility of  $\bar{R}$ . Next we bound the two error terms in (12),

- the **approximation error**:  $\inf_{\bar{R} \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(\bar{R}) - \mathcal{L}(R^*)|$ ;
- the **statistical error**:  $\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(R) - \hat{\mathcal{L}}(R)|$ .

Then Theorem 4.1 follows after bounding these two error terms.

### A. The approximation error

**Lemma 8.1.** *Suppose that (A1)-(A3) hold and the network parameters satisfy (N1) and (N2). Then,*

$$\inf_{\bar{R} \in \mathbf{R}_{\mathcal{D}, \mathcal{W}, \mathcal{S}, \mathcal{B}}} |\mathcal{L}(\bar{R}) - \mathcal{L}(R^*)| \leq 320C_1L_1B_1\sqrt{pd}n^{-1/(p+2)} + o(1). \quad (13)$$

as  $n \rightarrow \infty$ .

*Proof.* By (3) and (6) and the definition of  $\mathcal{L}$ , we have

$$\inf_{\bar{R} \in \mathbf{R}_{\mathcal{D}, \mathcal{W}, \mathcal{S}, \mathcal{B}}} |\mathcal{L}(\bar{R}) - \mathcal{L}(R^*)| \leq |\mathbb{D}_f(\mu_{\bar{R}_{\bar{\theta}}(\mathbf{x})} \|\gamma_d)| + |\mathcal{V}[R^*(\mathbf{x}), \mathbf{y}] - \mathcal{V}[\bar{R}_{\bar{\theta}}(\mathbf{x}), \mathbf{y}]|, \quad (14)$$

where  $\bar{R}_{\bar{\theta}} \in \mathbf{R}_{\mathcal{D}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$  is specified in Lemma 8.3 below. We finish the proof by (16) in Lemma 8.4 and (17) in Lemma 8.5, which will be proved below.  $\square$

**Lemma 8.2.** *For any function  $f : [-B, B]^p \rightarrow \mathbb{R}$  with Lipschitz constant  $L$  there exist a ReLU network  $\bar{f}$  with depth  $\mathcal{O}(12\mathcal{H} + C_{1,p})$  and width  $\mathcal{O}(C_{2,p}\mathcal{W})$  such that*

$$\|f - \bar{f}\|_{L^\infty} \leq 19L\sqrt{p}B(\mathcal{H}\mathcal{W})^{-2/p},$$

where  $C_{1,p} = 14 + 2p$ ,  $C_{2,p} = 3^{p+3}$ .

*Proof.* This Lemma follows directly from Theorem 1.1 of Shen (2020)  $\square$

**Lemma 8.3.** *Suppose that (A1) and (A3) hold and the network parameters satisfy (N1). Then, There exist a  $\bar{R}_{\bar{\theta}} \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  with the network parameters satisfying (N1) such that*

$$\|\bar{R}_{\bar{\theta}} - R^*\|_{L^2(\mu_{\mathbf{x}})} \leq 19L_1B_1\sqrt{pd}n^{-\frac{1}{p+2}}. \quad (15)$$

*Proof.* Let  $R_i^*(x)$  be the  $i$ -th entry of  $R^*(x) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ . By the assumption on  $R^*$ , it is

easy to check that  $R_i^*(x)$  is Lipschitz continuous on  $[-B_1, B_1]^d$  with the Lipschitz constant  $L_1$ . By Lemma 8.2, there exists a ReLU network  $\bar{R}_{\bar{\theta}_i}$  with depth  $\mathcal{O}(\mathcal{H})$  and width  $\mathcal{O}(\mathcal{W})$  such that

$$\|R_i^* - \bar{R}_{\bar{\theta}_i}\|_{L^\infty} \leq 19L_1B_1\sqrt{p}(\mathcal{H}\mathcal{W})^{-2/p}.$$

Then

$$\begin{aligned} \|R_i^* - \bar{R}_{\bar{\theta}_i}\|_{L^2(\mu_{\mathbf{x}})} &= \left[ \int (R_i^*(x) - \bar{R}_{\bar{\theta}_i}(x))^2 f_X(x) dx \right]^{1/2} \\ &\leq \|R_i^* - \bar{R}_{\bar{\theta}_i}\|_{L^\infty} \int f_X(x) dx \\ &\leq 19L_1B_1\sqrt{p}(\mathcal{H}\mathcal{W})^{-2/p}. \end{aligned}$$

Define  $\bar{R}_{\bar{\theta}} = [\bar{R}_{\bar{\theta}_1}, \dots, \bar{R}_{\bar{\theta}_d}] \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ . The above three display implies

$$\|\bar{R}_{\bar{\theta}} - \tilde{R}^*\|_{L^2(\mu_{\mathbf{x}})} \leq 19L_1B_1\sqrt{pd}(\mathcal{H}\mathcal{W})^{-2/p} \leq 19L_1B_1\sqrt{pd}n^{-1/(p+2)},$$

where in the last inequality we use the the choice of  $\mathcal{H}$  and  $\mathcal{W}$  in (N1). □

**Lemma 8.4.** *Suppose that (A1) and (A3) hold and the network parameters satisfy (N1). Then,*

$$|\mathcal{V}[R^*(\mathbf{x}), \mathbf{y}] - \mathcal{V}[\bar{R}_{\bar{\theta}}(\mathbf{x}), \mathbf{y}]| \leq 320C_1L_1B_1\sqrt{pd}n^{-1/(p+2)}. \quad (16)$$

*Proof.* Recall that Székely et al. (2007)

$$\begin{aligned} \mathcal{V}[\mathbf{z}, \mathbf{y}] &= \mathbb{E} [ \|\mathbf{z}_1 - \mathbf{z}_2\| \|\mathbf{y}_1 - \mathbf{y}_2\| ] - 2\mathbb{E} [ \|\mathbf{z}_1 - \mathbf{z}_2\| \|\mathbf{y}_1 - \mathbf{y}_3\| ] \\ &\quad + \mathbb{E} [ \|\mathbf{z}_1 - \mathbf{z}_2\| ] \mathbb{E} [ \|\mathbf{y}_1 - \mathbf{y}_2\| ], \end{aligned}$$

where  $(\mathbf{z}_i, \mathbf{y}_i), i = 1, 2, 3$  are i.i.d. copies of  $(\mathbf{z}, \mathbf{y})$ . We have

$$\begin{aligned}
& |\mathcal{V}[R^*(\mathbf{x}), \mathbf{y}] - \mathcal{V}[\bar{R}_{\bar{\theta}}(\mathbf{x}), \mathbf{y}]| \\
& \leq |\mathbb{E} [(\|R^*(\mathbf{x}_1) - R^*(\mathbf{x}_2)\| - \|\bar{R}_{\bar{\theta}}(\mathbf{x}_1) - \bar{R}_{\bar{\theta}}(\mathbf{x}_2)\|)|\mathbf{y}_1 - \mathbf{y}_2]| \\
& + 2|\mathbb{E} [(\|R^*(\mathbf{x}_1) - R^*(\mathbf{x}_2)\| - \|\bar{R}_{\bar{\theta}}(\mathbf{x}_1) - \bar{R}_{\bar{\theta}}(\mathbf{x}_2)\|)|\mathbf{y}_1 - \mathbf{y}_3]| \\
& + |\mathbb{E} [\|R^*(\mathbf{x}_1) - R^*(\mathbf{x}_2)\| - \|\bar{R}_{\bar{\theta}}(\mathbf{x}_1) - \bar{R}_{\bar{\theta}}(\mathbf{x}_2)\|] \mathbb{E} [\|\mathbf{y}_1 - \mathbf{y}_2\|]| \\
& \leq 8C_1 \mathbb{E} [|\|R^*(\mathbf{x}_1) - R^*(\mathbf{x}_2)\| - \|\bar{R}_{\bar{\theta}}(\mathbf{x}_1) - \bar{R}_{\bar{\theta}}(\mathbf{x}_2)\||] \\
& \leq 16C_1 \mathbb{E} [|\|R^*(\mathbf{x}) - \bar{R}_{\bar{\theta}}(\mathbf{x})\||] \\
& \leq 320C_1 L_1 B_1 \sqrt{pd} n^{-1/(p+2)}
\end{aligned}$$

where in the first and third inequalities we use the triangle inequality, and second one follows from the boundedness of  $\mathbf{y}$ , the last inequality is due to (15).  $\square$

**Lemma 8.5.** *Suppose that (A1) and (A2) hold and the network parameters satisfy (N1). Then,*

$$|\mathbb{D}_f(\mu_{\bar{R}_{\bar{\theta}}(\mathbf{x})} \|\gamma_d)| \rightarrow 0, \quad (17)$$

as  $n \rightarrow \infty$ .

*Proof.* By Lemma 8.3  $\bar{R}_{\bar{\theta}}$  can approximate  $R^*$  arbitrarily well as  $n \rightarrow \infty$ , the desired result follows from the fact that  $\mathbb{D}_f(\mu_{R^*(\mathbf{x})} \|\gamma_d) = 0$  and the continuity of  $\mathbb{D}_f(\mu_{R(\mathbf{x})} \|\gamma_d)$  on  $R$ . We present the sketch of the proof and omit the details here. Let  $r^*(z) = \frac{d\mu_{R^*(\mathbf{x})}}{d\gamma_d}(z)$  and  $\bar{r}(z) = \frac{d\mu_{\bar{R}_{\bar{\theta}}(\mathbf{x})}}{d\gamma_d}(z)$ . By definition we have

$$\mathbb{D}_f(\mu_{R^*(\mathbf{x})} \|\gamma_d) = \mathbb{E}_{W \sim \gamma_d}[f(r^*(W))]$$

We can represent  $\mathbb{D}_f(\mu_{\bar{R}_{\bar{\theta}}}||\gamma_d)$  similarly. Therefore,

$$\begin{aligned}
|\mathbb{D}_f(\mu_{\bar{R}_{\bar{\theta}(\mathbf{x})}}||\gamma_d)| &= |\mathbb{D}_f(\mu_{\bar{R}_{\bar{\theta}(\mathbf{x})}}||\gamma_d) - \mathbb{D}_f(\mu_{R^*(\mathbf{x})}||\gamma_d)| \\
&\leq \mathbb{E}_{W \sim \gamma_d} [|f(r^*(W)) - f(\bar{r}(W))|] \\
&\leq \int |f'(\tilde{r}(z))| |r^*(z) - \bar{r}(z)| d\gamma_d(z) \\
&\leq B_2 \int |r^*(z) - \bar{r}(z)| d\gamma_d(z),
\end{aligned}$$

where the second inequality we use mean value theorem and boundness assumption on  $f'(\tilde{r})$  in (A2). Then last inequality goes to zero due to continuity and the fact  $\bar{R}_{\bar{\theta}}$  converge to  $R^*$  as  $n \rightarrow \infty$  by Lemma 8.3.  $\square$

## B. The statistical error

**Lemma 8.6.** *Suppose that (A1)-(A2) hold and the network parameters satisfy (N1) and (N2). Then,*

$$\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(R) - \hat{\mathcal{L}}(R)| \leq C_{13}(2B_2 + B_3)n^{-\frac{1}{2+d}} + 19(1 + B_3)L_2\sqrt{d} \log nn^{-\frac{1}{d+2}} + 4C_6C_7C_{10}\mathcal{B}n^{-\frac{1}{p+2}} \quad (18)$$

*Proof.* By the definition and the triangle inequality we have

$$\begin{aligned}
\mathbb{E}[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathcal{L}(R) - \hat{\mathcal{L}}(R)|] &\leq \mathbb{E}[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\hat{\mathcal{V}}_n[R(\mathbf{x}), \mathbf{y}] - \mathcal{V}[(R(\mathbf{x}), \mathbf{y})]| \\
&\quad + \mathbb{E}[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\hat{\mathbb{D}}_f(\mu_{R(\mathbf{x})}||\gamma_d) - \mathbb{D}_f(\mu_{R(\mathbf{x})}||\gamma_d)|].
\end{aligned}$$

We finish the proof based on (19) in Lemma 8.7 and (24) in Lemma 8.8, which will be proved below.  $\square$

**Lemma 8.7.** *Suppose that (A1)-(A2) hold and the network parameters satisfy (N1) and (N2). Then,*

$$\mathbb{E}\left[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\widehat{\mathcal{V}}_n[R(\mathbf{x}), \mathbf{y}] - \mathcal{V}[R(\mathbf{x}), \mathbf{y}]|\right] \leq 4C_6 C_7 C_{10} \mathcal{B} n^{-\frac{1}{p+2}}. \quad (19)$$

*Proof.* We first fix some notation for simplicity. Denote  $O = (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^p \times \mathbb{R}^1$  and  $O_i = (\mathbf{x}_i, \mathbf{y}_i), i = 1, \dots, n$  are i.i.d copy of  $O$ , and denote  $\mu_{\mathbf{x}, \mathbf{y}}$  and  $\mathbb{P}^{\otimes n}$  as  $\mathbb{P}$  and  $\mathbb{P}^n$ , respectively.  $\forall R \in R \in \mathbf{R}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}$ , let  $\tilde{O} = (R(\mathbf{x}), \mathbf{y})$  and  $\tilde{O}_i = (R(\mathbf{x}_i), \mathbf{y}_i), i = 1, \dots, n$  are i.i.d copy of  $\tilde{O}$ . Define centered kernel  $\bar{h}_R : (\mathbb{R}^p \times \mathbb{R}^1)^{\otimes 4} \rightarrow \mathbb{R}$  as

$$\begin{aligned} \bar{h}_R(\tilde{O}_1, \tilde{O}_2, \tilde{O}_3, \tilde{O}_4) &= \frac{1}{4} \sum_{\substack{1 \leq i, j \leq 4, \\ i \neq j}} \|R(\mathbf{x}_i) - R(\mathbf{x}_j)\| |\mathbf{y}_i - \mathbf{y}_j| \\ &- \frac{1}{4} \sum_{i=1}^4 \left( \sum_{\substack{1 \leq j \leq 4, \\ j \neq i}} \|R(\mathbf{x}_i) - R(\mathbf{x}_j)\| \sum_{\substack{1 \leq j \leq 4, \\ i \neq j}} |\mathbf{y}_i - \mathbf{y}_j| \right) \\ &+ \frac{1}{24} \sum_{\substack{1 \leq i, j \leq 4, \\ i \neq j}} \|R(\mathbf{x}_i) - R(\mathbf{x}_j)\| \sum_{\substack{1 \leq i, j \leq 4, \\ i \neq j}} |\mathbf{y}_i - \mathbf{y}_j| - \mathcal{V}[R(\mathbf{x}), \mathbf{y}] \end{aligned} \quad (20)$$

Then, the centered  $U$ -statistics  $\widehat{\mathcal{V}}_n[R(\mathbf{x}), \mathbf{y}] - \mathcal{V}[R(\mathbf{x}), \mathbf{y}]$  can be represented as

$$\mathbb{U}_n(\bar{h}_R) = \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} \bar{h}_R(\tilde{O}_{i_1}, \tilde{O}_{i_2}, \tilde{O}_{i_3}, \tilde{O}_{i_4}).$$

Our goal is to bound the supremum of the centered  $U$ -process  $\mathbb{U}_n(\bar{h}_R)$  with the nondegenerate kernel  $\bar{h}_R$ . By the symmetrization randomization Theorem 3.5.3 in De la Pena and Giné (2012), we have

$$\mathbb{E}\left[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\mathbb{U}_n(\bar{h}_R)|\right] \leq C_5 \mathbb{E}\left[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} \left| \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} \epsilon_{i_1} \bar{h}_R(\tilde{O}_{i_1}, \tilde{O}_{i_2}, \tilde{O}_{i_3}, \tilde{O}_{i_4}) \right|\right], \quad (21)$$

where,  $\epsilon_{i_1}, i_1 = 1, \dots, n$  are i.i.d Rademacher variables that are also independent with  $\tilde{O}_i, i = 1, \dots, n$ . We finish the proof by upper bounding the above Rademacher process with the

metric entropy of  $R \in \mathbf{R}_{\bar{\mathcal{H}}, \bar{\mathcal{W}}, \bar{\mathcal{S}}}$ . To this end we need the following lemma.

**Lemma 8.8.** *If  $\xi_i, i = 1, \dots, m$  are  $m$  finite linear combinations of Rademacher variables  $\epsilon_j, j = 1, \dots, J$ . Then*

$$\mathbb{E}_{\epsilon_j, j=1, \dots, J} \max_{1 \leq i \leq m} |\xi_i| \leq C_6 (\log m)^{1/2} \max_{1 \leq i \leq m} (\mathbb{E} \xi_i^2)^{1/2}. \quad (22)$$

*Proof.* This result follows directly from Corollary 3.2.6 and inequality (4.3.1) in De la Pena and Giné (2012) with  $\Phi(x) = \exp(x^2)$ .  $\square$

By Lemma 8.3, we can assume the boundness of  $R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ , i.e., we can assume  $\|R\|_{L^\infty} \leq \mathcal{B} = 2\|R^*\|_{L^\infty}$  as  $n$  large enough. Then by the boundedness assumption on  $\mathbf{y}$ , we have that the kernel  $\bar{h}_R$  is also bounded, say

$$\|\bar{h}_R\|_{L^\infty} \leq C_7 \mathcal{B}. \quad (23)$$

$\forall R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ , define a random empirical measure (depends on  $O_i, i = 1, \dots, n$ )

$$e_{n,1}(R, \tilde{R}) = \mathbb{E}_{\epsilon_{i_1}, i_1=1, \dots, n} \left| \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} \epsilon_{i_1} (\bar{h}_R - \bar{h}_{\tilde{R}})(\tilde{O}_{i_1}, \dots, \tilde{O}_{i_4}) \right|.$$

Condition on  $O_i, i = 1, \dots, n$ , let  $\mathfrak{C}(\mathbf{R}, e_{n,1}, \delta)$  be the covering number of  $\mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  with respect to the empirical distance  $e_{n,1}$  at scale of  $\delta > 0$ . Denote  $\mathbf{R}_\delta$  as the covering set of

$\mathbf{R}_{\mathcal{D}, \mathcal{W}, \mathcal{S}}$  with cardinality of  $\mathfrak{C}(\mathbf{R}, e_{n,1}, \delta)$ . Then,

$$\begin{aligned}
& \mathbb{E}_{\epsilon_{i_1}} \left[ \sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} \left| \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} \epsilon_{i_1} \bar{h}_R(\tilde{O}_{i_1}, \tilde{O}_{i_2}, \tilde{O}_{i_3}, \tilde{O}_{i_4}) \right| \right] \\
& \leq \delta + \mathbb{E}_{\epsilon_{i_1}} \left[ \sup_{R \in \mathbf{R}_\delta} \left| \frac{1}{C_n^4} \sum_{1 \leq i_1 < i_2 < i_3 < i_4 \leq n} \epsilon_{i_1} \bar{h}_R(\tilde{O}_{i_1}, \tilde{O}_{i_2}, \tilde{O}_{i_3}, \tilde{O}_{i_4}) \right| \right] \\
& \leq \delta + C_6 \frac{1}{C_n^4} (\log \mathfrak{C}(\mathbf{R}, e_{n,1}, \delta))^{1/2} \max_{R \in \mathbf{R}_\delta} \left[ \sum_{i_1=1}^n \sum_{i_2 < i_3 < i_4} (\bar{h}_R(\tilde{O}_{i_1}, \tilde{O}_{i_2}, \tilde{O}_{i_3}, \tilde{O}_{i_4}))^2 \right]^{1/2} \\
& \leq \delta + C_6 C_7 \mathcal{B} (\log \mathfrak{C}(\mathbf{R}, e_{n,1}, \delta))^{1/2} \frac{1}{C_n^4} \left[ \frac{n(n!)^2}{((n-3)!)^2} \right]^{1/2} \\
& \leq \delta + 2C_6 C_7 \mathcal{B} (\log \mathfrak{C}(\mathbf{R}, e_{n,1}, \delta))^{1/2} / \sqrt{n} \\
& \leq \delta + 2C_6 C_7 \mathcal{B} (\text{VC}_{\mathbf{R}} \log \frac{2e\mathcal{B}n}{\delta \text{VC}_{\mathbf{R}}})^{1/2} / \sqrt{n} \\
& \leq \delta + C_6 C_7 C_{10} \mathcal{B} (\mathcal{H}\mathcal{S} \log \mathcal{S} \log \frac{\mathcal{B}n}{\delta \mathcal{D}\mathcal{S} \log \mathcal{S}})^{1/2} / \sqrt{n}.
\end{aligned}$$

where the first inequality follows from the triangle inequality, the second inequality uses (22), the third and fourth inequalities follow after some algebra, and the fifth inequality holds since  $\mathfrak{C}(\mathbf{R}, e_{n,1}, \delta) \leq \mathfrak{C}(\mathbf{R}, e_{n,\infty}, \delta)$  and the relationship between the metric entropy and the VC-dimension of the ReLU networks  $\mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  (Anthony and Bartlett, 2009), i.e.,

$$\log \mathfrak{C}(\mathbf{R}, e_{n,\infty}, \delta) \leq \text{VC}_{\mathbf{R}} \log \frac{2e\mathcal{B}n}{\delta \text{VC}_{\mathbf{R}}},$$

and the last inequality holds due to the upper bound of VC-dimension for the ReLU network  $\mathbf{R}_{\mathcal{D}, \mathcal{W}, \mathcal{S}}$  satisfying

$$C_8 \mathcal{H}\mathcal{S} \log \mathcal{S} \leq \text{VC}_{\mathbf{R}} \leq C_9 \mathcal{H}\mathcal{S} \log \mathcal{S},$$

see Bartlett et al. (2019). Then (19) holds by the selection of the network parameters and set  $\delta = \frac{1}{n}$  and some algebra.  $\square$

**Lemma 8.9.** *Suppose that (A1)-(A3) hold and the network parameters satisfy (N1) and (N2). Then,*

$$\mathbb{E}[\sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}} |\widehat{\mathbb{D}}_f(\mu_{R(\mathbf{x})} || \gamma_d) - \mathbb{D}_f(\mu_{R(\mathbf{x})} || \gamma_d)|] \leq C_{14}(L_2\sqrt{d} + B_2 + B_3)(n^{-\frac{2}{2+p}} + \log nn^{-\frac{2}{2+d}}) \quad (24)$$

*Proof.* For every  $R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ , let  $r(z) = \frac{d\mu_{R(\mathbf{x})}}{d\gamma_d}(z)$ ,  $g_R(z) = f'(r(z))$ . By assumption  $g_R(z) : \mathbb{R}^d \rightarrow \mathbb{R}$  is Lipschitz continuous with the Lipschitz constant  $L_2$  and  $\|g_R\|_{L^\infty} \leq B_2$ . By tail probability of Gaussian, we assume without loss of generality that  $\text{supp}(g_R) \subseteq [-\log n, \log n]^d$ . Then, by Lemma 8.2 there exists a  $\bar{D}_{\bar{\phi}} \in \mathbf{D}_{\bar{\mathcal{H}}, \bar{\mathcal{W}}, \bar{\mathcal{S}}}$  with the network parameters satisfying (N2) such that for  $\mathbf{z} \sim \gamma_d$  and  $\mathbf{z} \sim \mu_{R(\mathbf{x})}$ ,

$$\mathbb{E}_{\mathbf{z}}[|\bar{D}_{\bar{\phi}}(\mathbf{z}) - g_R(\mathbf{z})|] \leq 19L_2\sqrt{d}\log nn^{-\frac{1}{d+2}}. \quad (25)$$

By the above display, we can further assume that the element in  $\mathbf{D}_{\bar{\mathcal{H}}, \bar{\mathcal{W}}, \bar{\mathcal{S}}}$  is bounded by  $2B_2$  as  $n$  large enough. For any  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , define

$$\mathcal{E}(g) = \mathbb{E}_{\mathbf{x} \sim \mu_{\mathbf{x}}}[g(R(\mathbf{x}))] - \mathbb{E}_{W \sim \gamma_d}[f^*(g(W))],$$

$$\widehat{\mathcal{E}}(g) = \widehat{\mathcal{E}}(g, R) = \frac{1}{n} \sum_{i=1}^n [g(R(\mathbf{x}_i)) - f^*(g(W_i))].$$

By (6) we have

$$\mathcal{E}(g_R) = \mathbb{D}_f(\mu_{R(\mathbf{x})} || \gamma_d) = \sup_{\text{measurable } D: \mathbb{R}^d \rightarrow \mathbb{R}} \mathcal{E}(D). \quad (26)$$

Then,

$$\begin{aligned}
& |\mathbb{D}_f(\mu_{R(\mathbf{x})} || \gamma_d) - \widehat{\mathbb{D}}_f(\mu_{R(\mathbf{x})} || \gamma_d)| \\
&= |\mathcal{E}(g_R) - \max_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} \widehat{\mathcal{E}}(D_\phi)| \\
&\leq |\mathcal{E}(g_R) - \sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} \mathcal{E}(D_\phi)| + |\sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} \mathcal{E}(D_\phi) - \max_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} \widehat{\mathcal{E}}(D_\phi)| \\
&\leq |\mathcal{E}(g_R) - \mathcal{E}(\bar{D}_{\bar{\phi}})| + \sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} |\mathcal{E}(D_\phi) - \widehat{\mathcal{E}}(D_\phi)| \\
&\leq \mathbb{E}_{\mathbf{z} \sim \mu_{R(\mathbf{x})}} [|g_R - \bar{D}_{\bar{\phi}}|(\mathbf{z})] + \mathbb{E}_{W \sim \gamma_d} [|f^*(g_R) - f^*(\bar{D}_{\bar{\phi}})|(W)] + \sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{D}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}} |\mathcal{E}(D_\phi) - \widehat{\mathcal{E}}(D_\phi)| \\
&\leq 19(1 + B_3)L_2\sqrt{d} \log nn^{-\frac{1}{d+2}} + \sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}} |\mathcal{E}(D_\phi) - \widehat{\mathcal{E}}(D_\phi)|,
\end{aligned}$$

where we use the triangle inequality in the first inequality follows from the triangle inequality, the second inequality follows from  $\mathcal{E}(g_R) \geq \sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}} \mathcal{E}(D_\phi)$  due to (26) and the triangle inequality, the third inequality follows from the triangle inequality, and the last inequality follows from (25) and the mean value theorem.

We finish the proof by bounding the second term in probability in the last line above, i.e.,  $\sup_{D_\phi \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}} |\mathcal{E}(D_\phi) - \widehat{\mathcal{E}}(D_\phi)|$ . This can be done by bounding the empirical process

$$\mathbb{U}(D, R) = \mathbb{E} \left[ \sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}, \mathcal{B}}, D \in \mathcal{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}} |\mathcal{E}(D) - \widehat{\mathcal{E}}(D)| \right].$$

Let  $S = (\mathbf{x}, \mathbf{z}) \sim \mu_{\mathbf{x}} \otimes \gamma_d$  and  $S_i, i = 1, \dots, n$  be  $n$  i.i.d copy of  $S$ . Denote

$$b(D, R; S) = D(R(\mathbf{x})) - f^*(D(\mathbf{z})).$$

Then

$$\mathcal{E}(D, R) = \mathbb{E}_S[b(D, R; S)]$$

and

$$\widehat{\mathcal{E}}(D, R) = \frac{1}{n} \sum_{i=1}^n b(D, R; S_i).$$

Let

$$\mathcal{G}(\mathbf{D} \times \mathbf{R}) = \frac{1}{n} \mathbb{E}_{\{\epsilon_i, S_i\}_i^n} \left[ \sup_{R \in \mathbf{R}_{\mathcal{H}, \mathcal{W}, S, \mathcal{B}}, D \in \mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{S}, \tilde{\mathcal{B}}}} \left| \sum_{i=1}^n \epsilon_i b(D, R; S_i) \right| \right]$$

be the Rademacher complexity of  $\mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{S}} \times \mathbf{R}_{\mathcal{H}, \mathcal{W}, S}$  (Bartlett and Mendelson, 2002). Let  $\mathfrak{C}(\mathbf{D} \times \mathbf{R}, e_{n,1}, \delta)$  be the covering number of  $\mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{S}} \times \mathbf{R}_{\mathcal{H}, \mathcal{W}, S}$  with respect to the empirical distance (depends on  $S_i$ )

$$d_{n,1}((D, R), (\tilde{D}, \tilde{R})) = \frac{1}{n} \mathbb{E}_{\epsilon_i} \left[ \sum_{i=1}^n |\epsilon_i (b(D, R; S_i) - b(\tilde{D}, \tilde{R}; S_i))| \right]$$

at scale of  $\delta > 0$ . Let  $\mathbf{D}_\delta \times \mathbf{R}_\delta$  be such a converging set of  $\mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}} \times \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$ . Then,

$$\begin{aligned}
\mathbb{U}(D, R) &= 2\mathcal{G}(\mathbf{D} \times \mathbf{R}) \\
&= 2\mathbb{E}_{S_1, \dots, S_n} [\mathbb{E}_{\epsilon_i, i=1, \dots, n} [\mathcal{G}(\mathbf{R} \times \mathbf{D}) | (S_1, \dots, S_n)]] \\
&\leq 2\delta + \frac{2}{n} \mathbb{E}_{S_1, \dots, S_n} [\mathbb{E}_{\epsilon_i, i=1, \dots, n} [\sup_{(D, R) \in \mathbf{D}_\delta \times \mathbf{R}_\delta} |\sum_{i=1}^n \epsilon_i b(D, R; S_i)| | (S_1, \dots, S_n)]] \\
&\leq 2\delta + C_{12} \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} [(\log \mathfrak{C}(\mathbf{D} \times \mathbf{R}, e_{n,1}, \delta))^{1/2} \max_{(D, R) \in \mathbf{D}_\delta \times \mathbf{R}_\delta} [\sum_{i=1}^n b^2(D, R; S_i)]^{1/2}] \\
&\leq 2\delta + C_{12} \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} [(\log \mathfrak{C}(\mathbf{D} \times \mathbf{R}, e_{n,1}, \delta))^{1/2} \sqrt{n} (2B_2 + B_3)] \\
&\leq 2\delta + C_{12} \frac{1}{\sqrt{n}} (2B_2 + B_3) (\log \mathfrak{C}(\mathbf{D}, e_{n,1}, \delta) + \log \mathfrak{C}(\mathbf{R}, d_{n,1}, \delta))^{1/2} \\
&\leq 2\delta + C_{13} \frac{2B_2 + B_3}{\sqrt{n}} (\mathcal{H}\mathcal{S} \log \mathcal{S} \log \frac{\mathcal{B}n}{\delta \mathcal{H}\mathcal{S} \log \mathcal{S}} + \tilde{\mathcal{H}}\tilde{\mathcal{S}} \log \tilde{\mathcal{S}} \log \frac{2B_2n}{\delta \tilde{\mathcal{H}}\tilde{\mathcal{S}} \log \tilde{\mathcal{S}}})^{1/2}
\end{aligned}$$

where the first equality follows from the standard symmetrization technique, the first inequality holds due to the iteration law of conditional expectation, the second inequality follows from the triangle inequality, and the third inequality uses (22), the fourth inequality uses the fact that  $b(D, R; S)$  is bounded, i.e.,  $\|b(D, R; S)\|_{L^\infty} \leq 2B_2 + B_3$ , and the fifth inequality follows from some algebra, and the sixth inequality follows from  $\mathfrak{C}(\mathbf{R}, e_{n,1}, \delta) \leq \mathfrak{C}(\mathbf{R}, e_{n,\infty}, \delta)$  (similar result for  $\mathbf{D}$ ) and  $\log \mathfrak{C}(\mathbf{R}, e_{n,\infty}, \delta) \leq \text{VC}_{\mathbf{R}} \log \frac{2e\mathcal{B}n}{\delta \text{VC}_{\mathbf{R}}}$ , and  $\mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$  satisfying  $C_8 \mathcal{H}\mathcal{S} \log \mathcal{S} \leq \text{VC}_{\mathbf{R}} \leq C_9 \mathcal{H}\mathcal{S} \log \mathcal{S}$ , see Bartlett et al. (2019). Then (24) follows from the above display with the selection of the network parameters of  $\mathbf{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}}, \mathbf{R}_{\mathcal{H}, \mathcal{W}, \mathcal{S}}$  and with  $\delta = \frac{1}{n}$ .  $\square$

Finally, Theorem 4.1 is a direct consequence of (13) in Lemma 8.1 and (18) in Lemma 8.6. This completes the proof of Theorem 4.1.  $\square$

## References

- J. K. Aapo Hyvärinen, I. Erkki Oja, and P. J. Groenen. *Independent Component Analysis*. John-Wiley & Sons, Inc., New York, 2001.
- A. Achille and S. Soatto. Emergence of invariance and disentanglement in deep representations. *The Journal of Machine Learning Research*, 19(1):1947–1980, 2018.
- G. Alain and Y. Bengio. Understanding intermediate layers using linear classifier probes. *arXiv preprint arXiv:1610.01644*, 2016.
- S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B*, 28(1):131–142, 1966.
- B. Amos and J. Z. Kolter. A PyTorch Implementation of DenseNet. <https://github.com/bamos/densenet.pytorch>. Accessed: [20 Feb 2020].
- M. Anthony and P. L. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 2009.
- M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017.
- Y. Bai, T. Ma, and A. Risteski. Approximability of discriminators implies diversity in GANs. In *International Conference on Learning Representations*, 2019.
- P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482, 2002.
- P. L. Bartlett, N. Harvey, C. Liaw, and A. Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20:1–17, 2019.
- B. Bauer and M. Kohler. On deep learning as a remedy for the curse of dimensionality in nonparametric regression. *Ann. Statist.*, 47(4):2261–2285, 2019.
- Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- Y. Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Communications on Pure and Applied Mathematics*, 44(4):375–417, 1991.

- W. Bryc. *The Normal Distribution: Characterizations with Applications*. Lecture Notes in Statistics. Springer New York, 1995.
- M. Chen, W. Liao, H. Zha, and T. Zhao. Statistical guarantees of generative adversarial networks for distribution estimation. *arXiv preprint arXiv:2002.03938*, 2020.
- P. Comon. Independent component analysis, a new concept? *Signal Processing*, 36(3): 287–314, 1994.
- D. R. Cook. *Regression Graphics: Ideas for Studying Regressions Through Graphics*. Wiley Series in Probability and Statistics. Wiley, 1998. ISBN 9780471193654.
- D. R. Cook. Fisher lecture: Dimension reduction in regression. *Statistical Science*, 22(1): 1–26, 2007.
- D. R. Cook. Principal components, sufficient dimension reduction, and envelopes. *Annual Review of Statistics and Its Application*, 5:533–559, 2018.
- D. R. Cook and S. Weisberg. Sliced inverse regression for dimension reduction: comment. *Journal of the American Statistical Association*, 86(414):328–332, 1991.
- V. De la Pena and E. Giné. *Decoupling: from Dependence to Independence*. Springer Science & Business Media, 2012.
- M. H. Farrell, T. Liang, and S. Misra. Deep neural networks for estimation and inference. *Econometrica*, 89(1):181–213, 2021.
- R. A. Fisher. On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London, A*, 222:309–368, 1922.
- K. Fukumizu, A. Gretton, X. Sun, and B. Schölkopf. Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems*, pages 489–496, 2008.
- K. Fukumizu, F. R. Bach, and M. I. Jordan. Kernel dimension reduction in regression. *The Annals of Statistics*, 37(4):1871–1905, 2009.
- Y. Gao, Y. Jiao, Y. Wang, Y. Wang, C. Yang, and S. Zhang. Deep generative learning via variational gradient flow. In *International Conference on Machine Learning*, pages 2093–2101, 2019.

- Y. Gao, J. Huang, Y. Jiao, and J. Liu. Learning implicit generative models with theoretical guarantees. *arXiv preprint arXiv:2002.02862*, 2020.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27, pages 2672–2680. 2014.
- G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017a.
- G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4700–4708, 2017b.
- X. Huo and G. J. Székely. Fast computing for distance covariance. *Technometrics*, 58(4): 435–447, 2016.
- Y. Jiao, G. Shen, Y. Lin, and J. Huang. Deep nonparametric regression on approximately low-dimensional manifolds. *arXiv 2104.06708*, 2021.
- C. Jutten and J. Herault. Blind separation of sources, part i: An adaptive algorithm based on neuromimetic architecture. *Signal Processing*, 24(1):1–10, 1991.
- H. Kabir, M. Abdar, S. M. J. Jalali, A. Khosravi, A. F. Atiya, S. Nahavandi, and D. Srinivasan. Spinalnet: Deep neural network with gradual input. *arXiv preprint arXiv:2007.03347*, 2020.
- A. Keziou. Dual representation of  $\varphi$ -divergences and applications. *Comptes Rendus Mathématique*, 336(10):857–862, 2003.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling. Semi-supervised learning with deep generative models. In *Advances in neural information processing systems*, pages 3581–3589, 2014.
- A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Technical Report TR-2009, University of Toronto, Toronto.*, 2009.

- A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1:541–551, 1989.
- Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist>, 7:23, 2010.
- Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- K.-Y. Lee, B. Li, and F. Chiaromonte. A general theory for non-linear sufficient dimension reduction: formulation and estimation. *The Annals of Statistics*, 41(1):221–249, 2013.
- B. Li. *Sufficient Dimension Reduction: Methods and Applications with R*. CRC Press, 2018.
- B. Li, H. Zha, and F. Chiaromonte. Contour regression: a general approach to dimension reduction. *The Annals of Statistics*, 33(4):1580–1616, 2005.
- K.-C. Li. Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414):316–327, 1991.
- K.-C. Li. On principal hessian directions for data visualization and dimension reduction: Another application of stein’s lemma. *Journal of the American Statistical Association*, 87(420):1025–1039, 1992.
- Y. Li, K. Swersky, and R. Zemel. Generative moment matching networks. In *ICML*, 2015.
- T. Liang. How well generative adversarial networks learn distributions. 2020.
- Y. Ma and L. Zhu. Semiparametric approach to dimension reduction. *Journal of the American Statistical Association*, 107(497):168–179, 2012.
- Y. Ma and L. Zhu. Efficient estimation in sufficient dimension reduction. *Annals of Statistics*, 41(4):250–268, 2013a.
- Y. Ma and L. Zhu. A review on dimension reduction. *International Statistical Review*, 81(1):134–150, 2013b.

- J. C. Maxwell. Illustrations of the dynamical theory of gases. part i. on the motions and collisions of perfectly elastic spheres. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 19(124):19–32, 1860.
- R. J. McCann. Existence and uniqueness of monotone measure-preserving maps. *Duke Mathematical Journal*, 80(2):309–324, 1995.
- R. Mutihac and M. M. V. Hulle. A comparative survey on adaptive neural network algorithms for independent component analysis. *Romanian Reports in Physics*, 55:43–67, 2003.
- X. Nguyen, M. J. Wainwright, and M. I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861, 2010.
- K. Nordhausen and H. Oja. Independent component analysis: A statistical perspective. *WIREs Computational Statistics*, 10(5):e1440, 2018.
- S. Nowozin, B. Cseke, and R. Tomioka. f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, pages 271–279, 2016.
- S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2009.
- G. Philippis. *Regularity of optimal transport maps and applications*, volume 17. Springer Science & Business Media, 2013.
- S. Roberts and R. Everson. *Independent Component Analysis: Principles and Practice*. Cambridge University Press, 2001.
- T. R. Rockafellar. *Convex analysis*. Princeton University Press, 1970.
- A. Samarov and A. Tsybakov. Nonparametric independent component analysis. *Bernoulli*, 10(4):565 – 582, 2004.
- R. J. Samworth and M. Yuan. Independent component analysis via nonparametric maximum likelihood estimation. *The Annals of Statistics*, 40(6):2973 – 3002, 2012.
- J. Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61: 85–117, 2015.

- J. Schmidt-Hieber. Nonparametric regression using deep neural networks with ReLU activation function. *Ann. Statist.*, 48(4):1916–1921, 2020.
- Z. Shen. Deep network approximation characterized by number of neurons. *Communications in Computational Physics*, 28(5):1768–1811, 2020.
- D. J. Sutherland, H.-Y. Tung, H. Strathmann, S. De, A. Ramdas, A. Smola, and A. Gretton. Generative models and model criticism via optimized maximum mean discrepancy. In *ICLR*, 2017.
- T. Suzuki and M. Sugiyama. Sufficient dimension reduction via squared-loss mutual information estimation. *Neural computation*, 25(3):725–758, 2013.
- G. J. Székely and M. L. Rizzo. Brownian distance covariance. *The Annals of Applied Statistics*, 3(4):1236–1265, 2009.
- G. J. Székely, M. L. Rizzo, and N. K. Bakirov. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769–2794, 2007.
- C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu. A survey on deep transfer learning. In *International conference on artificial neural networks*, pages 270–279. Springer, 2018.
- L. Torrey and J. Shavlik. Transfer learning. In *Handbook of research on machine learning applications and trends: algorithms, methods, and techniques*, pages 242–264. IGI global, 2010.
- P. Vepakomma, C. Tonde, and A. Elgammal. Supervised dimensionality reduction via distance correlation maximization. *Electronic Journal of Statistics*, 12(1):960–984, 2018.
- C. Villani. *Optimal Transport: Old and New*, volume 338. Springer, 2008.
- R. Wang, A.-H. Karimi, and A. Ghodsi. Distance correlation autoencoder. In *2018 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2018.
- Y. Xia, H. Tong, W. Li, and L.-X. Zhu. An adaptive estimation of dimension reduction space. *Journal of the Royal Statistical Society Series B*, 64(3):363–410, 2002.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

- X. Yin and D. R. Cook. Dimension reduction for the conditional  $k$  th moment in regression. *Journal of the Royal Statistical Society: Series B*, 64(2):159–175, 2002.
- S. Zagoruyko and N. Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- P. Zhang, Q. Liu, D. Zhou, T. Xu, and X. He. On the discrimination-generalization tradeoff in GANs. In *International Conference on Learning Representations*, 2018.
- S. Zhao, J. Song, and S. Ermon. Infovae: Balancing learning and inference in variational autoencoders. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 5885–5892, 2019.
- L. Zhu, L. Zhu, and Z. Feng. Dimension reduction in regressions through cumulative slicing estimation. *Journal of the American Statistical Association*, 105(492):1455–1466, 2010.