

A Deep Generative Approach to Conditional Sampling

Xingyu Zhou^{*†} Yuling Jiao^{*‡} Jin Liu[§] and Jian Huang[¶]

October 22, 2021

Abstract

We propose a deep generative approach to sampling from a conditional distribution based on a unified formulation of conditional distribution and generalized nonparametric regression function using the noise-outsourcing lemma. The proposed approach aims at learning a conditional generator so that a random sample from the target conditional distribution can be obtained by the action of the conditional generator on a sample drawn from a reference distribution. The conditional generator is estimated nonparametrically with neural networks by matching appropriate joint distributions using Kullback-Liebler divergence. An appealing aspect of our method is that it allows either of or both the predictor and the response to be high-dimensional and can handle both continuous and discrete type predictors and responses. We show that the proposed method is consistent in the sense that the conditional generator converges in distribution to the underlying conditional distribution under mild conditions. Our numerical experiments with simulated and benchmark image data validate the proposed method and demonstrate that it outperforms several existing conditional density estimation methods.

Keywords: Distribution matching; Generative learning; High-dimensional data; Nonparametric estimation; Neural networks.

^{*}Xingyu Zhou and Yuling Jiao contributed equally to this work.

[†]Department of Statistics and Actuarial Science, University of Iowa, Iowa City, Iowa 52242, USA. Email: xingyu-zhou@uiowa.edu

[‡]School of Mathematics and Statistics, Wuhan University, Wuhan, Hubei, China 430072. Email: yulingjiaomath@whu.edu.cn

[§]Centre for Quantitative Medicine, Duke-NUS Medical School, Singapore. Email: jin.liu@duke-nus.edu.sg

[¶]Corresponding author. Department of Statistics and Actuarial Science, University of Iowa, Iowa City, Iowa 52242, USA. Email: jian-huang@uiowa.edu

1 Introduction

A fundamental problem in statistics and machine learning is how to model the relationship between a response Y and a predictor X . Such a model can be used for predicting the values of Y based on the new observations of X and for assessing the variation in Y for a given value of X . Regression models that focus on estimating the conditional mean or median of the response given the predictor have been widely used for such purposes. However, in problems when the conditional distribution is multimodal or asymmetrical, conditional mean and median are no longer adequate for modeling the relationship between Y and X . In general, to completely understand how the response depends on the predictor, it becomes necessary to learn the conditional distribution, which provides a full description of the relationship between the response variable and the predictor. Conditional distribution also plays a central role in many important areas, including representation learning (Bengio et al., 2013), sufficient dimension reduction (Li, 1991; Cook, 1998), graphical models (Bishop, 2006), nonlinear independent component analysis (Hyvärinen and Pajunen, 1999), among others.

In this paper, we propose a nonparametric generative approach to sampling from a conditional distribution. For convenience, we shall refer to the proposed method as the generative conditional distribution sampler (GCDS). For a given value of the predictor $X = x$, GCDS aims at estimating a function $G(\eta, x)$ of η and x , where η is a random variable from a simple reference distribution such as normal or uniform, such that $G(\eta, x)$ follows the conditional distribution of Y given $X = x$. Such a function is called a conditional generator. To sample from the conditional distribution, we only need to calculate $G(\eta, x)$ after generating η from the reference distribution. Therefore, the conditional generator G contains all the information about the conditional distribution of Y given X . We estimate

the conditional generator G nonparametrically using neural networks. An appealing feature of GCDS is that it is applicable to the settings when either of or both X and Y are high-dimensional such as in the problems of image data analysis.

There is an extensive literature on nonparametric conditional density estimation. The prevailing approaches are based on smoothing methods, including kernel smoothing and local polynomials (Rosenblatt, 1969; Scott, 1992; Hyndman et al., 1996; Chen and Linton, 2001; Hall and Yao, 2005; Bott and Kohler, 2017). Typically, the joint density of (X, Y) and the marginal density of X are first estimated using unconditional kernel density estimators. Then, the conditional density estimator is obtained as the ratio of the estimated joint density over the estimated marginal density. Another approach is to transform the problem of estimating a conditional density to a suitably formulated regression problem (Fan et al., 1996; Fan and Yim, 2004) and use the method for nonparametric regression for conditional density estimation. Nearest neighbors method has also been used in estimating conditional density and conditional quantiles through kernel smoothing (Zhao and Liu, 1985; Bhattacharya and Gangopadhyay, 1990). Approaches based on expanding the conditional density function in terms of certain basis functions have also been developed (Sugiyama et al., 2010; Izbicki and Lee, 2016). The method proposed by Izbicki et al. (2017) approximates the conditional density using orthogonal basis functions and transform the problem of conditional density estimation into a regression problem. A common feature of these methods is that they seek to estimate the functional form of the conditional density.

However, these existing conditional density estimation methods do not work well for problems with high-dimensional data. In particular, they suffer from the so-called “curse of dimensionality”, that is, their performance deteriorates dramatically as the dimensionality of the dependent variable or the response variable becomes relatively large. Indeed,

most conditional density estimators can only effectively handle up to a few covariates. In addition, most of these methods focus on the case when the response Y is a scalar variable and cannot handle the case of a high-dimensional response vector.

The proposed approach is inspired by the recently developed generative adversarial networks (GAN) (Goodfellow et al., 2014). Instead of estimating the functional form of the conditional density, GCDS is a generative learning approach that seeks to estimate a conditional sampler. The basis of GCDS is a unified formulation of the conditional density estimation and the generalized nonparametric regression based on the noise-outsourcing lemma in probability theory (Kallenberg, 2002; Austin, 2015). By this lemma, the problem of nonparametric conditional density estimation is equivalent to a generalized nonparametric regression problem. This equivalency implies that, for any given $X = x$, we can estimate the conditional generator $G(\eta, x)$ so that a random sample from the conditional distribution can be obtained based on this function using a random sample η from a reference distribution, such as the uniform or the standard normal distribution. The estimation of the conditional generator G is achieved through matching appropriate joint distributions using the Kullback-Liebler divergence and its variational form. We take advantage of the abilities of neural networks in approximating high-dimensional functions and estimate G nonparametrically using deep neural networks.

There are several advantages of GCDS over the classical methods for conditional density estimation. First, there is no restriction on the dimensionality of the response variable, while the classical methods typically only consider the case of a scalar response variable. Indeed, our methods allow either of or both the predictor and the response to be high-dimensional. Second, GCDS can handle both continuous and discrete type predictors and responses. Third, since our method learns a generative function for the underlying condi-

tional distribution based on a simple reference distribution, it is easy to obtain estimates of the summary measures of the underlying conditional distribution, including the conditional moments and quantiles by Monte Carlo. In comparison, it is cumbersome to do so based on the traditional conditional density estimation methods, since it involves numerical integrations that are difficult to implement in high-dimensional settings. Fourth, we demonstrate that the proposed method works for complex and high-dimensional data problems such as image generation and reconstruction. The traditional conditional density estimation methods are not able to deal with such problems. Finally, we show that GCDS is consistent in the sense that the samples it generates converge weakly to the underlying target conditional distribution. To the best of our knowledge, such a result is the first of its kind in the context of deep generative learning.

In the remainder of this paper, we first describe a generative representation of conditional distribution based on the noise-outsourcing lemma, and explain that sampling from a conditional distribution can be achieved by using a conditional generator. This provides the theoretical foundation for the distribution matching method proposed in Section 3. The distribution matching is carried out by using the variational form of the f -divergence, which includes the Kullback-Liebler divergence as a special case. In Section 4 we establish the consistency of the conditional generator in the sense that the joint distribution of X and generated sample converges to the joint distribution of (X, Y) . In Section 5 we conduct extensive simulation studies to evaluate the finite sample performance of the proposed method and illustrate its application to an image generation and reconstruction problems using benchmark image data. Concluding remarks are given in Section 6. Additional numerical experiment results and technical details are provided in the supplementary material.

2 Generative representation of conditional distribution

Consider a pair of random vectors $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, where X is a vector of predictors and Y is a vector of response variables or labels. For regression problems, we have $\mathcal{Y} \subseteq \mathbb{R}^q$ with $q \geq 1$; for classification problems, \mathcal{Y} is a set of finite many labels. We assume the space $\mathcal{X} \subseteq \mathbb{R}^d$ with $d \geq 1$. The predictor X can contain both continuous and categorical components. Suppose $(X, Y) \sim P_{X,Y}$ with marginal distributions $X \sim P_X$ and $Y \sim P_Y$. Denote the conditional distribution of Y given X by $P_{Y|X}$. For a given value x of X , we also write the conditional distribution as $P_{Y|X=x}$. Let η be a random vector independent of X with a known distribution P_η . For example, we can take P_η to be the standard multivariate normal $N(\mathbf{0}, \mathbf{I}_m)$, for a given $m \geq 1$. We note that m does not need to be the same as q , the dimension of Y .

Our goal is to find a function $G : \mathbb{R}^m \times \mathcal{X} \mapsto \mathcal{Y}$ such that the conditional distribution of $G(\eta, X)$ given $X = x$ is the same as the conditional distribution of Y given $X = x$. Since η is independent of X , it is equivalent to finding a G such that

$$G(\eta, x) \sim P_{Y|X=x}, \quad x \in \mathcal{X}, \quad (1)$$

Thus to sample from the conditional distribution $P_{Y|X=x}$, we can first sample an $\eta \sim P_\eta$, then calculate $G(\eta, x)$. The resulting value $G(\eta, x)$ is a sample from $P_{Y|X=x}$.

Does such a function G exist? The existence of G is guaranteed by the noise-outsourcing lemma from probability theory under minimal conditions (Theorem 5.10 in Kallenberg (2002), Lemma 3.1 in Austin (2015)).

Lemma 2.1. (*Noise-outsourcing lemma*). *Let (X, Y) be a random pair taking values in $\mathcal{X} \times \mathcal{Y}$ with joint distribution $P_{X,Y}$. Suppose \mathcal{Y} is a standard Borel space. Then there exist*

a random vector $\eta \sim P_\eta = N(\mathbf{0}, \mathbf{I}_m)$ for any given $m \geq 1$ and a Borel-measurable function $G : \mathbb{R}^m \times \mathcal{X} \rightarrow \mathcal{Y}$ such that η is independent of X and

$$(X, Y) = (X, G(\eta, X)) \text{ almost surely.} \quad (2)$$

Because η and X are independent, any G satisfying (2) also satisfies (1), that is, $G(\cdot, x)$ is a conditional generator for $P_{Y|X=x}$. In the original noise-outsourcing lemma, the distribution P_η is a uniform distribution on $[0, 1]$. In Lemma 2.1, P_η is taken to be $N(\mathbf{0}, \mathbf{I}_m)$ with $m \geq 1$. This is more convenient in applying GCDS when it is better to use a random vector instead of a random scalar as the noise source. In the supplementary material, we show that we can indeed take P_η to be $N(\mathbf{0}, \mathbf{I}_m)$ based on the original noise-outsourcing lemma.

Lemma 2.1 provides a unified view of conditional distribution estimation and (generalized) nonparametric regression. To see this, it is informative to reverse the order of (1) and write it as

$$Y|X = x \sim G(\eta, x), \quad x \in \mathcal{X}. \quad (3)$$

This expression shows that the problem of finding G is similar to that of estimating a generalized regression function nonparametrically by matching the conditional distributions. Therefore, G can also be considered a generalized regression function. The standard formulation of nonparametric regression with an additive error is a special case of (3). Indeed, if we assume $G(\eta, x) = G_0(x) + \eta$ with $\mathbb{E}(\eta|X) = 0$, then (3) leads to the standard nonparametric regression model $\mathbb{E}(Y|X = x) = G_0(x)$.

Sampling from a conditional distribution generally cannot be done by simply using the existing methods for sampling from an unconditional distribution. This can be illustrated as follows. For any given x , the problem is to find a function $G_x(\eta)$ such that $G_x(\eta) \sim P_{Y|X=x}$,

where we use x as the subscript of G to indicate that the form of G depends on x . If X is a discrete random variable and only takes the values in a finite set, then we can simply find the function G_x for each x using the existing generative methods such as GAN (Goodfellow et al., 2014). However, this is not feasible if X is a continuous-type random variable. In general, the methods for generating samples from an unconditional distribution cannot be directly applied to find a function of η for generating samples from the conditional distribution $P_{Y|X}$.

To get around this difficulty, we note that matching the conditional distribution of $G(\eta, x)$ with $P_{Y|x}$ for a given $x \in \mathcal{X}$ is equivalent to matching the joint distribution of $(X, G(\eta, X))$ and the joint distribution of (X, Y) , if the same marginal distribution of X is involved. This can be easily seen as follows. Let $T = G(\eta, X)$. Then $P_{T|X} = P_{Y|X}$ if and only if $P_{T|X}P_X = P_{Y|X}P_X$ on the support of (X, Y) , that is, $P_{X,T} = P_{X,Y}$. We summarize this simple but key observation in the following lemma.

Lemma 2.2. *Suppose that η is independent of X . Then $G(\eta, x) \sim P_{Y|X=x}$, $x \in \mathcal{X}$ if and only if*

$$(X, G(\eta, X)) \sim (X, Y). \quad (4)$$

Because of (4), we refer to G as a conditional generator, since given $X = x$, $G(\eta, x) \sim P_{Y|X=x}$. Lemma 2.2 shows that finding a G such that (1) holds amounts to finding a G such that the joint distribution of $(X, G(\eta, X))$ is the same as that of (X, Y) .

It is clear that the conditional generator satisfying (4) contains all the information about the conditional distribution of Y given X . For example, consider the conditional expectation $g(x) = \mathbb{E}(Y|X = x)$ and the conditional variance $v(x) = \text{Var}(Y|X = x)$. By (4), we have $g(x) = \mathbb{E}_{\eta \sim P_\eta} G(\eta, x)$ and $v(x) = \text{Var}_\eta[G(x, \eta)]$. Therefore, we can calculate g and v based on G . Although it is difficult to calculate these functions exactly, it is easy

to approximate them via Monte Carlo. Specifically, let η_1, \dots, η_J be a random sample generated from P_η , then we can approximate $g(x)$ and $v(x)$ by

$$\tilde{g}(x) = \frac{1}{J} \sum_{j=1}^J G(\eta_j, x) \quad \text{and} \quad \tilde{v}(x) = \frac{1}{J} \sum_{j=1}^J [G(\eta_j, x) - \tilde{g}(x)]^2. \quad (5)$$

Since it is easy and inexpensive to generate random samples from P_η , for any given x we can easily accurately approximate the summary measures such as moments and quantiles of the conditional distribution $P_{Y|X=x}$ based on $\{G(\eta_j, x), j = 1, \dots, J\}$ for a sufficiently large J .

3 Distribution matching estimation

3.1 Adversarial generative networks

The generative adversarial networks (GAN) (Goodfellow et al., 2014) is an approach to learning a high-dimensional (unconditional) distribution. It is formulated as a minimax adversarial game between two players, a generator G and a discriminator D . The discriminator D is parameterized using a neural network that serves as a witness to distinguish between a sample Y from the data distribution and a sample from the generative model. The generator $G(\eta)$ maps samples η from the reference distribution P_η to the data distribution. The generator G is trained to maximally confuse the discriminator into believing that samples it generates come from the data distribution. Formally, GAN solves the minimax optimization problem:

$$\min_G \max_D \mathbb{E}_{Y \sim P_{\text{data}}} \log D(Y) + \mathbb{E}_{\eta \sim P_\eta} \log[1 - D(G(\eta))]. \quad (6)$$

Conditional generative adversarial networks (cGAN) (Mirza and Osindero, 2014) estimate the distribution of the images conditioning on some auxiliary information, especially class labels. Similar to GAN, it solves a two-player minimax game using an objective function with the same form as (6). See equation (2) in Mirza and Osindero (2014). cGAN performs the conditioning by feeding the class label information into the neural networks for the discriminator and the generator as additional input layer. However, cGAN does not work well for data generation with continuous conditions.

3.2 f -divergence and its variational form

While the minimax formulation has an attractive intuitive interpretation as a two-player game between the generator and the discriminator, it is helpful to understand it as the dual form of the primal problem of minimizing the Jensen-Shannon divergence between the data distribution and the distribution of the generator (Goodfellow et al., 2014). By considering a general discrepancy measure between two distributions such as the f -divergence, one can formulate a class of generative learning methods including GAN as a special case (Nowozin et al., 2016).

By Lemma 2.2, we can estimate G by matching the distribution of $(X, G(\eta, X))$ with the distribution of (X, Y) . For this purpose, we first describe the f -divergence and its variational form. Let P and Q be two probability distributions on \mathbb{R}^d . Let p and q be the density functions of P and Q with respect to a common dominant measure, respectively. Suppose Q is absolutely continuous with respect to P . The f -divergence (Ali and Silvey, 1966) of Q with respect to P is defined by

$$\mathbb{D}_f(q||p) = \int f\left(\frac{q(z)}{p(z)}\right) p(z) dz, \tag{7}$$

where $f : \mathbb{R}^+ \rightarrow \mathbb{R}$ is a convex function with $f(1) = 0$ and is strictly convex at $x = 1$. A basic property of the f -divergence following from Jensen's inequality is that $\mathbb{D}_f(q||p) \geq 0$ for every q, p and $\mathbb{D}_f(q||p) = 0$ if and only if $q = p$.

The Kullbak-Liebler (KL) divergence is an important special case with $f(x) = x \log x$, which has a simple expression

$$\mathbb{D}_{\text{KL}}(q||p) = \int \frac{q(z)}{p(z)} \log \left(\frac{q(z)}{p(z)} \right) p(z) dz = \int \log \left(\frac{q(z)}{p(z)} \right) q(z) dz. \quad (8)$$

Let $r = q/p$ be the density ratio of the densities q and p . It is convenient to express the KL divergence as

$$\mathbb{D}_{\text{KL}}(q||p) = \int \log \left(\frac{q(z)}{p(z)} \right) q(z) dz = \mathbb{E}_{Z \sim q}[\log r(Z)].$$

A useful representation of the f -divergence is its variational form. We will use it to construct an objective function for training the conditional generator G . The variational form is based on the Fenchel conjugate of f (Rockafellar, 1970), defined as $f^*(t) = \sup_{x \in \mathbb{R}} \{tx - f(x)\}, t \in \mathbb{R}$. Then the f -divergence has the following variational formulation (Keziou, 2003; Nguyen et al., 2010).

Lemma 3.1. *Let \mathcal{D} be a class of measurable functions $D : \mathbb{R}^d \rightarrow \mathbb{R}$. Suppose f is a differentiable convex function. Then*

$$\mathbb{D}_f(q||p) \geq \sup_{D \in \mathcal{D}} [\mathbb{E}_{Z \sim q} D(Z) - \mathbb{E}_{W \sim p} f^*(D(W))], \quad (9)$$

where the equality holds if and only if $f'(q/p) \in \mathcal{D}$ and the supremum is attained at $D^* = f'(q/p)$.

Commonly used divergence measures, including the KL divergence, the Jensen-Shannon (JS) divergence and the χ^2 -divergence, can be considered special cases of f -divergence. We give a proof of Lemma 3.1 and the expressions of the conjugate functions and variational forms of these divergence measures in the supplementary material.

3.3 Distribution matching estimation via f -divergence

We now apply Lemmas 2.2 and 3.1 to construct the objective function for estimating the conditional generator G . Let $p_{X,G(\eta,X)}$ and $p_{X,Y}$ be the densities of $(X, G(\eta, X))$ and (X, Y) , respectively. At the population level, we seek to find a conditional generator G^* that minimizes the f -divergence $\mathbb{D}_f(p_{X,G(\eta,X)} \| p_{X,Y})$.

Lemma 3.2. *A function $G^* : \mathbb{R}^m \times \mathcal{X} \rightarrow \mathcal{Y}$ is a minimizer of the f -divergence $\mathbb{D}_f(p_{X,G(\eta,X)} \| p_{X,Y})$,*

$$G^* \in \underset{G}{\operatorname{argmin}} \mathbb{D}_f(p_{X,G(\eta,X)} \| p_{X,Y}) \quad (10)$$

if and only if $p_{X,G^(\eta,X)} = p_{X,Y}$, that is, $(X, G^*(\eta, X)) \sim (X, Y)$.*

This lemma is a direct consequence from the properties of the f -divergence. Let

$$r(z) = \frac{p_{X,G(\eta,X)}(z)}{p_{X,Y}(z)}, \quad z \in \mathbb{R}^d \times \mathbb{R}^q. \quad (11)$$

be the density ratio of $p_{X,G(\eta,X)}$ over $p_{X,Y}$. We only focus on the KL divergence below. By (8), we have

$$\mathbb{D}_{\text{KL}}(p_{X,G(\eta,X)} \| p_{X,Y}) = \mathbb{E}_{(X,\eta) \sim p_X p_\eta} [\log r(X, G(\eta, X))].$$

Our goal is to minimize an empirical version of $\mathbb{D}_{\text{KL}}(p_{X,G(\eta,X)} \| p_{X,Y})$ with respect to G . The minimizer will serve as an estimator of G . The KL divergence depends on the unknown

density functions $p_{X,G(\eta,X)}$ and $p_{X,Y}$ only through the density ratio r or the log-density ratio. Denote the log-density ratio by $D = \log r$. To estimate G , we will also need to estimate D . We note that estimating density ratio is usually easier than estimating individual densities separately. The log-density ratio D can be intuitively interpreted as a discriminator that quantifies the difference between the distributions of $(X, G(\eta, X))$ and (X, Y) .

Therefore, in our problem the loss function determined by the log-density ratio D needs to be estimated along with the parameter of interest G . For this purpose, we consider the variational form of the KL divergence. The dual of $f(x) = x \log x$ is $f^*(t) = \exp(t - 1)$ (Nguyen et al., 2010). By Lemma 3.1, we can write the variational representation of the KL-divergence as

$$\begin{aligned} \mathbb{D}_{\text{KL}}(p_{X,G(\eta,X)} \| p_{X,Y}) &= \sup_D \{ \mathbb{E}_{(X,\eta) \sim P_X P_\eta} [D(X, G(\eta, X))] - \mathbb{E}_{(X,Y) \sim P_{X,Y}} [\exp(D(X, Y) - 1)] \} \\ &= \sup_D \{ \mathbb{E}_{(X,\eta) \sim P_X P_\eta} [D(X, G(\eta, X))] - \mathbb{E}_{(X,Y) \sim P_{X,Y}} [\exp(D(X, Y))] \} + 1, \end{aligned} \quad (12)$$

where the second equality follows by change of variables from $D - 1$ to D in the supremum operation. For the purpose of estimating G by minimizing the KL divergence, we can ignore the constant 1 in (12). So we consider the criterion

$$\mathcal{L}(G, D) = \mathbb{E}_{(X,\eta) \sim P_X P_\eta} [D(X, G(\eta, X))] - \mathbb{E}_{(X,Y) \sim P_{X,Y}} [\exp(D(X, Y))]. \quad (13)$$

The variational form is convenient since it is easy to obtain its empirical version when random samples are available. Then at the population level, the target conditional generator

G^* and the target discriminator D^* are characterized by the minimax problem

$$(G^*, D^*) = \underset{G}{\operatorname{argmin}} \underset{D}{\operatorname{argmin}} \mathcal{L}(G, D). \quad (14)$$

Suppose that $\{(X_i, Y_i), i = 1, \dots, n\}$ are i.i.d. $P_{X,Y}$ and $\{\eta_i, i = 1, \dots, n\}$ are independently generated from P_η . We consider the following empirical version of $\mathcal{L}(G, D)$:

$$\widehat{\mathcal{L}}(G, D) = \frac{1}{n} \sum_{i=1}^n D(X_i, G(\eta_i, X_i)) - \frac{1}{n} \sum_{i=1}^n \exp(D(X_i, Y_i)). \quad (15)$$

We estimate G nonparametrically using feedforward neural networks (FNN) (Schmidhuber, 2015) based on the objective function $\widehat{\mathcal{L}}(G, D)$ in (15). We use two FNNs: the conditional generator network G_θ with parameter θ for estimating G and a second network D_ϕ with parameter ϕ for estimating the discriminator 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.

- The generator network G_θ : let $\mathcal{G} \equiv \mathcal{G}_{\mathcal{H}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$ be the set of ReLU neural networks $G_\theta : \mathbb{R}^m \times \mathbb{R}^d \rightarrow \mathbb{R}^q$ with parameter θ , depth \mathcal{H} , width \mathcal{W} , size \mathcal{S} , and $\|G_\theta\|_\infty \leq \mathcal{B}$. 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 $w_0 = d$ is the dimension of the input data and $w_{\mathcal{H}} = q$ 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. For multilayer perceptrons with equal-width hidden layers except the output layer, we have $\mathcal{S} = \mathcal{W}(m + 1) + (\mathcal{W}^2 + \mathcal{W})(\mathcal{H} - 1) + \mathcal{W} + q$.
- The discriminator network D_ϕ : Similarly, denote $\mathcal{D} \equiv \mathcal{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}$ as the set of ReLU

neural networks $D_\phi : \mathbb{R}^d \times \mathbb{R}^q \rightarrow \mathbb{R}$, with parameter ϕ , depth $\tilde{\mathcal{H}}$, width $\tilde{\mathcal{W}}$, size $\tilde{\mathcal{S}}$, and $\|D_\phi\|_\infty \leq \tilde{\mathcal{B}}$.

Then θ and ϕ are estimated by solving the empirical version of the minimax problem (14), that is,

$$(\hat{\theta}, \hat{\phi}) = \underset{\theta}{\operatorname{argmin}} \underset{\phi}{\operatorname{argmax}} \hat{\mathcal{L}}(G_\theta, D_\phi). \quad (16)$$

The estimated conditional generator is $\hat{G} = G_{\hat{\theta}}$ and the estimated discriminator is $\hat{D} = D_{\hat{\phi}}$. It is natural to compute $(\hat{\theta}, \hat{\phi})$ by alternately minimizing $L(\theta, \phi)$ with respect to θ fixing ϕ and maximizing $L(\theta, \phi)$ with respect to ϕ fixing θ . We provide the implementation details in Section 5.

4 Weak convergence of conditional sampler

In this section, we provide sufficient conditions under which $(X, \hat{G}_n(\eta, X))$ converges in distribution to (X, Y) . This implies that for given $X = x$ with $p_X(x) > 0$, the conditional distribution of $\hat{G}_n(\eta, x)$ given $X = x$ converges to the conditional distribution of Y given $X = x$. We focus on the case when X and Y are continuous-type random vectors. We establish a slightly stronger result by showing that the total variation norm

$$\|p_{X, \hat{G}(\eta, X)} - p_{X, Y}\|_{L_1} = \int_{\mathcal{X} \times \mathcal{Y}} |p_{X, \hat{G}(\eta, X)}(x, y) - p_{X, Y}(x, y)| dx dy$$

converges to zero.

Let $\mathcal{L}(G, D)$ be defined in (13). For any measurable function $G : \mathbb{R}^m \times \mathbb{R}^d \mapsto \mathbb{R}^q$, define

$$\mathbb{L}(G) = \sup_D \mathcal{L}(G, D). \quad (17)$$

For a fixed G , let p_{XG} be the joint density of $(X, G(\eta, X))$. Lemma 3.1 implies that the optimal D is $D^*(z) = \log(p_{XG}(z)/p_{XY}(z)) = \log r(z)$. Thus the optimal discriminator is the log-likelihood ratio serving as a critic of the resemblance between p_{XY} and p_{XG} . Substituting this expression into (17) yields $\mathbb{L}(G) = \mathbb{E}_{(X,\eta) \sim P_X P_\eta}[\log r(X, G(\eta, X))]$. Let $G^* \in \operatorname{argmin}_G \mathbb{L}(G)$. We have $P_{(X,G^*(X,\eta))} = P_{X,Y}$ by Lemmas 3.1 and 3.2.

We assume the following conditions.

- (A1) The target conditional generator $G^* : \mathbb{R}^m \times \mathcal{X} \rightarrow \mathcal{Y}$ is continuous with $\|G^*\|_\infty \leq C_0$ for some constant $0 < C_0 < \infty$.
- (A2) For any $G \in \mathcal{G} \equiv \mathcal{G}_{\mathcal{D}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$, $r_G(z) = p_{X,G(\eta,X)}(z)/p_{X,Y}(z) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is continuous and $0 < C_1 \leq r_G(z) \leq C_2$ for some constants $0 < C_1 \leq C_2 < \infty$.

We also make the following assumptions on the network parameters of the conditional generator G_θ and the discriminator D_ϕ .

- (N1) The network parameters of \mathcal{G} satisfies

$$\mathcal{H}\mathcal{W} \rightarrow \infty \quad \text{and} \quad \frac{\mathcal{B}\mathcal{S}\mathcal{H} \log(\mathcal{S}) \log n}{n} \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

- (N2) The network parameters of \mathcal{D} satisfies

$$\tilde{\mathcal{H}}\tilde{\mathcal{W}} \rightarrow \infty \quad \text{and} \quad \frac{\tilde{\mathcal{B}}\tilde{\mathcal{S}}\tilde{\mathcal{H}} \log(\tilde{\mathcal{S}}) \log n}{n} \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Theorem 4.1. *Suppose that the assumptions (A1) and (A2) hold. If the network parameters of \mathcal{G} and \mathcal{D} satisfies the specifications (N1) and (N2), then*

$$\mathbb{E}_{(X_i, Y_i, \eta_i)_{i=1}^n} \|p_{X, \hat{G}_\theta(\eta, X)} - p_{X,Y}\|_{L_1}^2 \rightarrow 0, \quad \text{as } n \rightarrow \infty. \quad (18)$$

A direct corollary of Theorem 4.1 is the following convergence result in terms of the conditional density functions.

Corollary 4.1. *Suppose that the assumptions (A1) and (A2) hold and the network parameters of \mathcal{G} and \mathcal{D} satisfies the specifications (N1) and (N2). Then*

$$\mathbb{E}_{X \sim P_X} \left[\int_{\mathcal{Y}} \left| p_{\hat{G}_{\theta}(\eta, X)}(y|X) - p_{Y|X}(y|X) \right| dy \right] \rightarrow_P 0, \text{ as } n \rightarrow \infty.$$

Theorem 4.1 and Corollary 4.1 provide strong theoretical support for the proposed method under mild conditions. They are proved using the empirical process method (Bartlett and Mendelson, 2002; Bartlett et al., 2019) and the recent results on approximating smooth functions by deep neural networks (Shen et al., 2020). The main challenge of the proof is that the objective function is a minimax process indexed by two classes of neural networks. Details are given in the supplementary material.

Conditions (A1) and (A2) are mild regularity conditions that are often assumed in non-parametric estimation problems. Conditions (N1) and (N2) concern the depths, widths and sizes of the generator and the discriminator networks. For the generator These conditions require that the size of the network increases with the sample size, the product of the depth and the width increases with the sample size. We note that the conditions are flexible with respect to the network architecture. In particular, they allow either the depth or the width remain fixed. For example, we can have a deep network with fixed width or a wide network with fixed depth. A restriction of the conditions is that they require the network size to be smaller than the sample size. This restriction stems from the use of empirical process theory (Van der Vaart and Wellner, 1996; Bartlett et al., 2019; Bartlett and Mendelson, 2002) to control the stochastic error of the estimated generator and discriminator.

In nonparametric regression, there has been much recent work on convergence analysis of nonparametric estimators using deep neural networks. Two types of assumptions on the underlying model have been used in the analysis. The first type of assumptions postulates that the regression function has a compositional structure so that the intrinsic dimension of the function is lower than the ambient dimension (Bauer and Kohler, 2019; Kohler and Langer, 2020; Schmidt-Hieber, 2020; Shen et al., 2021a). The second type assumes that the distribution of X is supported on a lower-dimensional manifold (Chen et al., 2019; Nakada and Imaizumi, 2019; Jiao et al., 2021; Shen et al., 2021b). These works also require the size of the neural network used in the nonparametric regression to be smaller than the sample size to ensure the consistency of the estimators.

In the current problem of sampling from conditional distributions, convergence analysis is substantially more difficult than that in nonparametric regression. The main reason is that in the current setting, optimization is a minimax problem that leads to an *estimated nonparametric loss function*, i.e., there is a second neural network involved for estimating the discriminator in the dual form of KL divergence, in addition to the neural network for estimating the conditional generator. In comparison, nonparametric regression with a given loss function such as least squares loss only has a single neural network for estimating the regression function.

We now give a high-level description of the proof for Theorem 4.1, the details are provided in the supplementary material. Lemma 3.2 implies that $\mathbb{L}(G^*) = 0$. For notational simplicity, write $\hat{G} = \hat{G}_\theta$. By Pinsker’s inequality (Tsybakov, 2008), we have

$$\|p_{X, \hat{G}(\eta, X)} - p_{X, Y}\|_{L^1}^2 \leq 2(\mathbb{L}(\hat{G}) - \mathbb{L}(G^*)). \quad (19)$$

So it suffices to show that the right side in (19) converges to zero in expectation. By the

definition of $\mathbb{L}(G)$, we can write the excess risk as

$$\mathbb{L}(\hat{G}) - \mathbb{L}(G^*) = \sup_D \mathcal{L}(\hat{G}, D) - \sup_D \mathcal{L}(G^*, D).$$

A key step in the proof is the following decomposition of the excess risk. For any $\bar{G} \in \mathcal{G}$, we decompose the right side in (19) as:

$$\begin{aligned} \mathbb{L}(\hat{G}) - \mathbb{L}(G^*) &= \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) \\ &\quad + \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\hat{G}, D) \end{aligned} \tag{20}$$

$$+ \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\bar{G}, D) \tag{21}$$

$$+ \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\bar{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\bar{G}, D) \tag{22}$$

$$+ \sup_{D \in \mathcal{D}} \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(\bar{G}, D) \tag{23}$$

$$+ \sup_D \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(G^*, D).$$

Since the terms in (21) and (23) are nonpositive, and the terms in (20) and (22) are smaller than $\sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)|$, we have

$$\begin{aligned} \mathbb{L}(\hat{G}) - \mathbb{L}(G^*) &\leq \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) + 2 \sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)| \\ &\quad + \sup_D \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(G^*, D). \end{aligned}$$

Note that \bar{G} is arbitrary. By taking infimum with respect to \bar{G} over \mathcal{G} on both sides of the above display, we obtain

$$\mathbb{L}(\hat{G}) - \mathbb{L}(G^*) \leq \Delta_1 + \Delta_2 + \Delta_3, \tag{24}$$

where

$$\begin{aligned}\Delta_1 &= \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D), \\ \Delta_2 &= 2 \sup_{G \in \mathcal{G}, D \in \mathcal{D}} |\mathcal{L}(G, D) - \hat{\mathcal{L}}(G, D)|, \\ \Delta_3 &= \inf_{\tilde{G} \in \tilde{\mathcal{G}}} [\mathbb{L}(\tilde{G}) - \mathbb{L}(G^*)].\end{aligned}$$

The first and the third terms Δ_1 and Δ_3 are the approximation errors and the second term Δ_2 is the statistical error. In the supplementary material, we derive (24) and show that these error terms converge to zero.

5 Implementation

We describe the implementation of GCDS. For training the generator G_θ and the discriminator D_ϕ , we use the rectified linear unit (ReLU) as the activation function in G_θ and D_ϕ . We train the discriminator and the generator iteratively by updating θ and ϕ alternately as follows:

- (a) Fix θ , update the discriminator by ascending the stochastic gradient of the loss (15) with respect to ϕ .
- (b) Fix ϕ , update the generator by descending the stochastic gradient of the loss (15) with respect to θ .

The training process is described below.

Algorithm: Training GCDS

Input: (a) Pairs $\{(X_i, Y_i), i = 1, \dots, n\}$; (b) Samples $\{\eta\}_{i=1}^n$ from P_η

Output: Conditional generator $G_{\hat{\theta}}$ and discriminator $D_{\hat{\theta}}$

While not converged do

- Compute $\tilde{Y}_i = G_{\theta}(\eta_i, X_i)$, $i = 1, 2, \dots, n$. Let $S_1 = \{(X_i, Z_i, V_i) = (X_i, Y_i, 1), i = 1, \dots, n\}$ and $S_2 = \{(X_i, Z_i, V_i) = (X_{i-n}, \tilde{Y}_i, -1), i = n + 1, \dots, 2n\}$.
- Randomly select $B/2$ samples from S_1 and another $B/2$ samples from S_2 . Denote the subscripts of the selected samples by $\{b_i : i = 1, \dots, B\}$.
- Update D_{ϕ} by ascending its stochastic gradient:

$$\nabla_{\phi} \left\{ \frac{1}{B} \sum_{i=1}^B [D_{\phi}(X_{b_i}, Z_{b_i}) \mathbb{1}_{\{V_{b_i}=-1\}} - \exp(D_{\phi}(X_{b_i}, Z_{b_i})) \mathbb{1}_{\{V_{b_i}=1\}}] \right\}.$$

- Randomly select B samples from $\{(X_i, Y_i), i = 1, \dots, n\}$. Denote the subscripts of the selected samples by $\{b_i : i = 1, \dots, B\}$
- Update G_{θ} by descending its stochastic gradient:

$$\nabla_{\theta} \left\{ \frac{1}{B} \sum_{i=1}^B D_{\phi}(X_{b_i}, G_{\theta}(\eta_{b_i}, X_{b_i})) \right\}.$$

End while

We implement the GCDS algorithm in TensorFlow (Abadi et al., 2016).

6 Numerical experiments

In this section, we carry out numerical experiments to assess the performance of GCDS. We use both simulated and real datasets in the experiments. In addition to the results reported in this section, additional numerical results are given in the online supplementary material,

including results from experiments evaluating how the performance of GCDS depends on the network architecture, the dimension m of the noise vector η and the sample size n .

6.1 Simulation studies

We conduct simulation studies to evaluate the finite sample performance of GCDS. We also compare it with several existing conditional density estimation methods, including the nearest neighbor kernel conditional density estimation (NNKCDE, Dalmaso et al. (2020)), the conditional kernel density estimation (CKDE, implemented in the R package `np`, Hall et al. (2004)), and the basis expansion method FlexCode (Izbicki et al., 2017)). We implement GCDS in TensorFlow (Abadi et al., 2016) and use the stochastic gradient descent algorithm Adam (Kingma and Ba, 2015) in training the neural networks. We use the conditional distributions based on the following models in the simulation studies.

1. (M1). A nonlinear model with an additive error term:

$$Y = X_1^2 + \exp(X_2 + X_3/3) + \sin(X_4 + X_5) + \varepsilon, \varepsilon \sim N(0, 1).$$

2. (M2). A model with an additive error term whose variance depends on the predictors:

$$Y = X_1^2 + \exp((X_2 + X_3/3)) + X_4 - X_5 + (0.5 + X_2^2/2 + X_5^2/2) \times \varepsilon, \varepsilon \sim N(0, 1).$$

3. (M3). A model with a multiplicative non-Gaussian error term:

$$Y = (5 + X_1^2/3 + X_2^2 + X_3^2 + X_4 + X_5) * \exp(0.5 \times \varepsilon), \text{ where} \\ \varepsilon \sim \mathbb{I}_{\{U < 0.5\}} \times N(-2, 1) + \mathbb{I}_{\{U > 0.5\}} \times N(2, 1) \text{ with } U \sim \text{Uniform}(0, 1), X \in \mathbb{R}^{30}.$$

4. (M4). A mixture of two normal distributions:

$$Y = \mathbb{I}_{\{U < 0.5\}} N(-X_1, 0.25^2) + \mathbb{I}_{\{U > 0.5\}} N(X_1, 0.25^2), \text{ where } U \sim \text{Uniform}(0, 1).$$

In each of the models above, the covariate vector X is generated from standard multivariate normal distribution.

The neural networks used in the simulations are specified as follows. For models (M1)-(M3), the generator network has 1 hidden layer with width 50, and the discriminator has 2 hidden layers with widths (50, 25); for models (M4), the generator network has 2 hidden layers with widths (40, 15), and the discriminator has 2 hidden layers with widths (50, 25). For the values of m , the dimension of the noise random vector η , we set $m = 3$ for models (M1)-(M3), and $m = 4$ for model (M4).

For the conditional density estimation method NNKCDE, the tuning parameters are chosen using cross-validation. The bandwidth of the conditional kernel density estimator CKDE is chosen by the rule-of-thumb using the standard formula $h_j = 1.06\sigma_j n^{-1/(2*K+J)}$ where σ_j is a measure of spread of the j th continuous variable defined as $\min(SD, IQR/1.349)$, n the number of observations, K the order of the kernel, and J the number of continuous variables. The basis expansion based method FlexCode uses Fourier basis. The maximum number of bases is 40 and the actual number of bases is selected using cross-validation.

We calculate the mean squared error (MSE) of the estimated conditional mean $E(Y|X)$ and the estimated conditional standard deviation $SD(Y|X)$. We use a test data set $\{x_1, \dots, x_k\}$ of size $k = 2000$. The MSE of the estimated conditional mean is $MSE(\text{mean}) = (1/k) \sum_{i=1}^k [\hat{E}(Y|X = x_i) - E(Y|X = x_i)]^2$. For GCDS, the estimate of $E(Y|X = x)$ is based on (5) using Monte Carlo. For other methods, the estimate is calculated by numerical integration $\hat{E}(Y|x) = \int yf(y|x)dy$ using 1000 subdivisions. Similarly, the MSE of the estimated conditional standard deviation is $MSE(\text{sd}) = (1/k) \sum_{i=1}^k [\hat{SD}(Y|X = x_i) - SD(Y|X = x_i)]^2$.

		GCDS	NNKCDE	CKDE	FlexCode
M1	Mean	0.259 (0.015)	1.367(0.010)	0.491(0.024)	0.610(0.008)
	SD	0.022 (0.004)	0.258(0.004)	0.233(0.005)	0.170(0.007)
M2	Mean	0.312 (0.017)	4.668(0.046)	1.707(0.060)	2.408(0.063)
	SD	0.247 (0.012)	0.793(0.008)	0.857(0.017)	2.384(0.602)
M3	Mean	3.377 (0.196)	4.926(0.080)	39.084(0.929)	9.015(0.341)
	SD	2.082 (0.126)	8.131(0.235)	15.70(0.488)	11.53(1.140)
M4	Mean	0.016(0.003)	0.004 (0.001)	0.063(0.002)	0.006(0.002)
	SD	0.027 (0.005)	0.131(0.001)	0.076(0.001)	0.046(0.001)

Table 1: Mean squared error(MSE) of the estimated conditional mean, the estimated standard deviation and the corresponding simulation standard errors (in parentheses). The smallest MSEs are in bold font.

For GCDS, we first generate J samples $\{\eta_j : j = 1, \dots, J\}$ from the reference distribution P_η and calculate conditional samples $\{\hat{G}(\eta_j, x_i), j = 1, \dots, J\}$. We take $J = 10,000$. The estimated conditional standard deviation is calculated as the sample standard deviation of the conditional samples. The estimated conditional standard deviation of other methods are computed by numerical integration $\hat{SD}(Y|x_i) = \sqrt{\int [y - \hat{E}(Y|x_i)]^2 f(y|x_i) dy}$ using 1000 subdivisions.

We repeat the simulations 10 times. The average MSEs and simulation standard errors are summarized in Table 1. We see that, comparing with CKDE and FlexCode and NNKCED, GCDS has the smallest MSEs for estimating conditional mean and conditional SD in most cases.

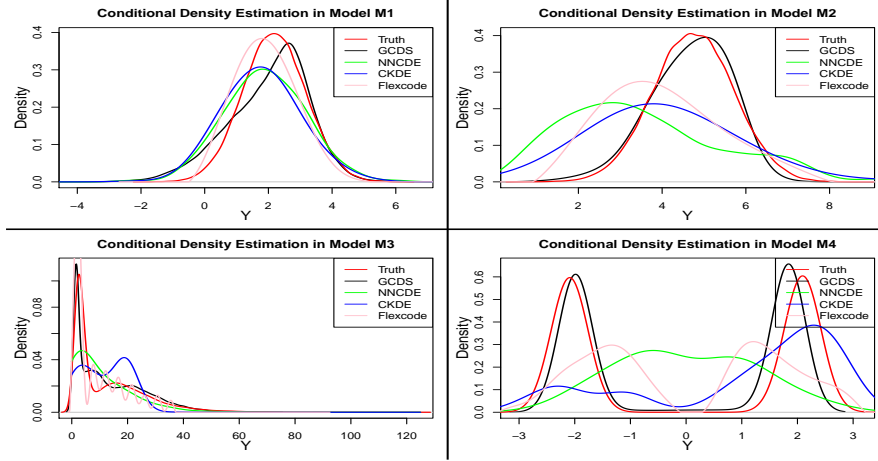


Figure 1: Comparison of density estimation in models (M1) to (M4). The conditional density function corresponding to GCDS is estimated using the samples generated from GCDS with kernel smoothing.

In Figure 1, we display the estimated conditional density functions for a randomly generated value of X . The true conditional distributions of $Y|X$ for models (M1) to (M4) are: (M1), $N(2.19, 1)$; (M2), $N(4.75, 0.96^2)$; (M3): the mixture of half $7.42 \times \log\text{-normal}(-1, 0.5^2)$ and half $7.42 \times \log\text{-normal}(1, 0.5^2)$; (M4): the mixture of half $N(-2.09, 0.25^2)$ and half $N(2.09, 0.25^2)$. The conditional density function corresponding to GCDS is estimated based the samples generated based on GCDS using kernel smoothing. This plot shows that GCDS yields better conditional density estimates than CKDE, NNCDE and FlexCode.

6.2 The abalone dataset

The abalone dataset is available at UCI machine learning repository (Dua and Graff, 2017). It contains the number of rings of abalone and other physical measurements. The age of

abalone is determined by cutting the shell through the cone, staining it, and counting the number of rings through a microscope, a time-consuming process. Other measurements, which are easier to obtain, are used to predict the number of rings that determines the age. This dataset contains 9 variables. They are *sex*, *length*, *diameter*, *height*, *whole weight*, *shucked weight*, *viscera weight*, *shell weight* and *rings*. Except for the categorical variable *sex*, all the other variables are continuous. The variable *sex* codes three groups: female, male and infant, since the gender of an infant abalone is not known. The sample size is 4177. We use 90% of the data for training and 10% of the data as the testing set. The neural networks used in the analysis are specified as follows: the generator network is a fully connected network with 2 hidden layers with widths 50 and 20; the discriminator networks is a fully connected network with 2 hidden layers with widths 50 and 25. The dimension of the noise vector of the noise vector is set to be $m = 5$.

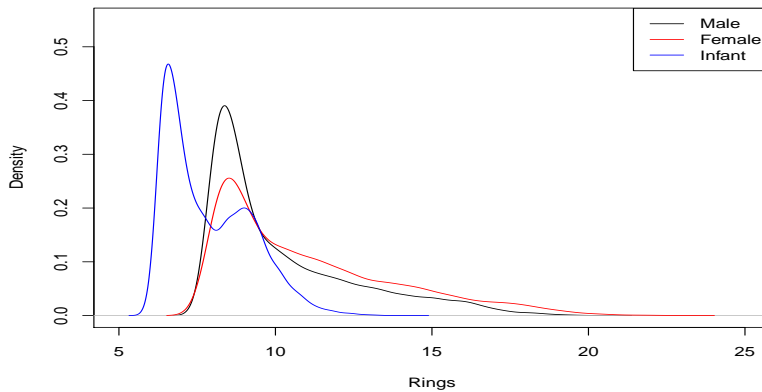


Figure 2: Estimated conditional densities for the female, male and infant groups in the abalone dataset. Each line represents the kernel conditional density estimation based on the samples generated using GCDS given the group average values of the covariates.

We take *rings* as the response $Y \in \mathbb{R}$ and the other measurements as the covariate vector $X \in \mathbb{R}^9$. Figure 2 shows the estimated conditional density based on the training

dataset for 3 groups: female, male and infant, at the value of the group means of the remaining covariates. We see that the values of *rings* of the infant group are smaller than those of the female and male groups. The female abalones tend to have slightly higher numbers of rings than male abalones. In addition, the conditional distributions are skewed to the right for all the three groups.

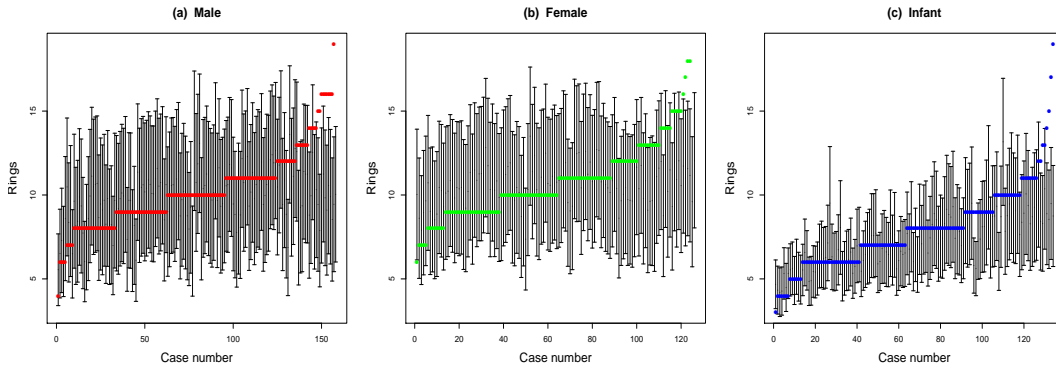


Figure 3: The prediction intervals for the testing set. All 418 abalones in the testing set are divided into three groups, (a) male, (b) female, and (c) infant.

To examine the prediction performance of the estimated conditional density, we construct the 90% prediction interval for the number of rings of each abalone in the testing set. The prediction intervals are shown in Figure 3. The actual number of rings are plotted as a solid dot. The actual coverage for all 418 cases in the testing set is 89.71%, close to the nominal level of 90%. The numbers of rings that are not covered by the prediction intervals are the largest ones in each group.

6.3 MNIST handwritten digits

We now illustrate the application of GCDS to high-dimensional data problems and demonstrate that it can easily handle the models when either of both of X and Y are high-dimensional. The data example we use is the MNIST handwritten digits dataset (LeCun et al., 2010), which contains 60,000 images for training and 10,000 images for testing. The images are stored in 28×28 matrices with gray color intensity from 0 to 1. Each image is paired with a label in $\{0, 1, \dots, 9\}$. We use GCDS to perform two tasks: generating images from labels and reconstructing the missing part of an image.

Generating images from labels We generate images of handwritten digits given the label. In this problem, the predictor X is a categorical variable representing the ten digits: $\{0, 1, \dots, 9\}$ and the response Y represents 28×28 images. We use one-hot vectors in \mathbb{R}^{10} to represent these ten categories. So the dimension of X is 10 and the dimension of Y is $28 \times 28 = 784$. The response $Y \in [0, 1]^{28 \times 28}$ is a matrix representing the intensity values. For the discriminator D , we use a convolutional neural network (CNN) with 3 convolution layers with 128, 256, and 256 filters to extract the features of the image and then concatenate with the label information (repeated 10 times to match the dimension of the features). The concatenated information is sent to a fully connected layer and then to the output layer. For the generator G , we concatenate the label information with random noise of dimension 100. Then it is fed to a CNN with 3 deconvolution layers with 256, 128, and 1 filters.

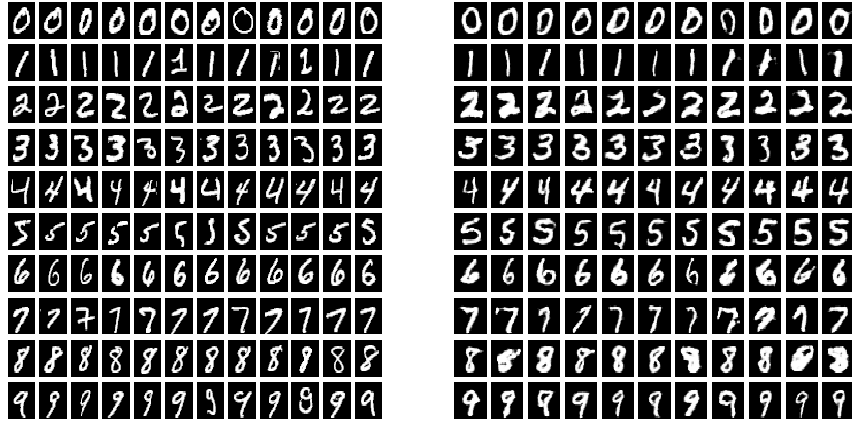


Figure 4: MNIST dataset: real images (left panel) and generated images given the labels (right panel).

Figure 4 shows the real images (left panel) and generated images (right panel). We see that the generated images are similar to the real images and it is hard to distinguish the generated ones from the real images. Also, there are some differences in the generated images, reflecting the random variations in the generating process.

Reconstructing missing part of an image We now illustrate using GCDS to reconstruct an image when part of the image is missing with the MNIST dataset. Suppose we only observe $1/4$, $1/2$ or $3/4$ of an image and would like to reconstruct the missing part of the image. For this problem, let X be the observed part of the image and let Y be the missing part of the image. Our goal is to reconstruct Y based on X . For the discriminator, we use two convolutional networks to process X and Y separately. The filters are then concatenated together and fed into another convolution layer and fully-connected layer before output. For the generator, X is processed by a fully-connected layer followed by 3 deconvolution layers.

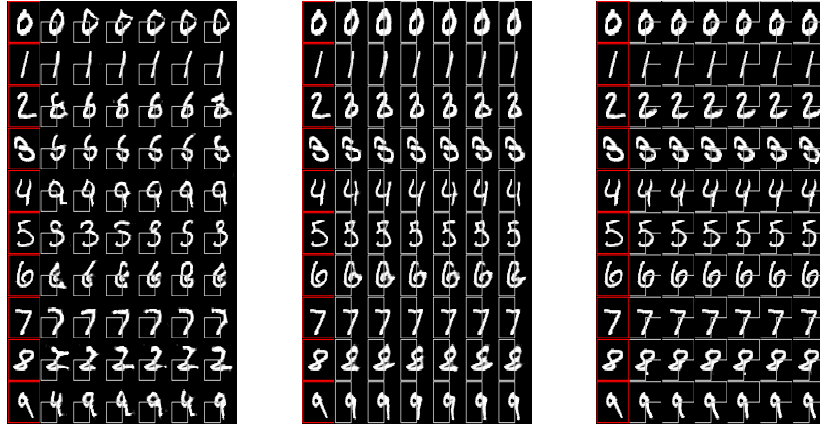


Figure 5: Reconstructed images given partial image in MNIST dataset. The first column in each panel consists of the true images, the other columns give the constructed images. In the left panel, the left lower $1/4$ of the image is given; in the middle panel, the left $1/2$ of the image is given; in the right panel, $3/4$ of the image is given.

In Figure 5, three plots from left to right corresponds to the situations when $1/4$, $1/2$ and $3/4$ of an image are given. In each subplot, the first column contains the true images in the testing set. The gray boxes show the given areas. Each row contains six reconstructions of the image. The digits “0”, “1” and “7” are easy to reconstruct. Even when only $1/4$ of their images are given, GCDS can correctly reconstruct them. The other digits are more difficult. If only $1/4$ of their images are given, it is impossible to reconstruct them. However, as the given area increases from $1/4$ to $1/2$ and then $3/4$ of the images, GCDS is able to reconstruct the images correctly, and the reconstructed images become less variable and more similar to the true image. For example, for the digit “2”, if only the left lower $1/4$ of the image is given, the reconstructed images tend to be incorrect; the reconstruction is only successful when $3/4$ of the image is given.

7 Conclusion

In this paper we propose GCDS, a generative approach to sampling from a conditional distribution. We provide theoretical support for GCDS by showing that the conditional generator converge in distribution to the underlying target conditional distribution under mild conditions. Our numerical experiments demonstrate that it works well in a variety of situations from the standard nonparametric conditional density estimation problems to more complex image data problems.

Several questions deserve further investigation. First, it would be interesting to derive the convergence rate of the sampling distribution that strengthens the consistency result in Theorem 4.1 and provide conditions under which the number of coefficients in the deep neural network is allowed to be greater than the sample size. Second, while the conditional generator provides all the information of the conditional distribution, it is still useful to obtain an estimate of the functional form of the conditional density. How to obtain a good estimator of the conditional density function in the present framework is an open question, especially when the dimension of (X, Y) is high. Finally, as a proof of concept we demonstrated that GCDS yields reasonable results for some simple image analysis tasks with the MNIST dataset. It would be interesting to apply GCDS to more complex image analysis problems. We intend to study these challenging problems in the future.

Supplementary material

Additional numerical experiment results and technical details are provided in the supplementary material.

Acknowledgements

The work of X. Zhou and J. Huang is supported in part by the U.S. National Science Foundation grant DMS-1916199. The work of Y. Jiao is supported by the National Science Foundation of China grants No.11871474 and No.61701547. The work of J. Liu is supported by the Duke-NUS Graduate Medical School WBS: R-913-200-098-263 and MOE2016-M2-2-029 from the Ministry of Education, Singapore.

Supplementary material

In this Supplementary Material, we provide additional numerical experiment results, present the variational forms of three commonly used f -divergences, explain the difference between the f -divergence definition used in this paper and that in Nguyen et al. (2010), and prove Lemma 2.1, Lemma 3.1, and Theorem 4.1.

A Additional numerical experiments

A.1 Additional simulation results

A.1.1 Conditional quantile estimation

We consider the conditional distributions based on the following models as given in (M1)-(M4) in the main text:

1. Model (M1):

$$Y = X_1^2 + \exp(X_2 + X_3/3) + \sin(X_4 + X_5) + \varepsilon, \quad \varepsilon \sim N(0, 1).$$

This is a nonlinear model with an additive error term.

2. Model (M2):

$$Y = X_1^2 + \exp((X_2 + X_3/3)) + X_4 - X_5 + (0.5 + X_2^2/2 + X_5^2/2) \times \varepsilon, \quad \varepsilon \sim N(0, 1).$$

This model has an additive error term whose variance depends on the predictors.

3. Model (M3):

$$Y = (5 + X_1^2/3 + X_2^2 + X_3^2 + X_4 + X_5) * \exp(0.5 \times \varepsilon)$$

$$\varepsilon \sim \mathbb{I}_{\{U < 0.5\}} \times N(-2, 1) + \mathbb{I}_{\{U > 0.5\}} \times N(2, 1) \text{ where } U \sim \text{Uniform}(0, 1),$$

with $p = 30$. This model has a multiplicative non-Gaussian error term.

4. Model (M4):

$$Y = \mathbb{I}_{\{U < 0.5\}} N(-X_1, 0.25^2) + \mathbb{I}_{\{U > 0.5\}} N(X_1, 0.25^2),$$

$$U \sim \text{Uniform}(0, 1).$$

In this example, the conditional distribution is a mixture of two normal distributions.

In each of the models above, the covariate vector X is generated from standard multivariate normal distribution.

In models (M1) to (M3), the conditional generator G is parameterized by a one-layer neural network; in model (M4), it is parameterized by a two-layer fully connected neural network. The log-density ratio function D is parameterized by a two-layer fully connected neural network. We use the stochastic gradient descent algorithm as implemented in Adam

in training the neural networks.

We now examine the estimation of the conditional quantiles for a given $X = x$ defined as $q_\tau = F_{Y|x}^{-1}(\tau|x)$ for $0 < \tau < 1$.

The conditional quantiles can be estimated easily using Monte Carlo. For a given $X = x$, we generate J a random sample of i.i.d. η_1, \dots, η_J from the reference distribution P_η . Then the estimated quantiles of the conditional distribution at $X = x$ is calculated based on $\hat{G}(x, \eta_j), j = 1, \dots, J$. We take $J = 10,000$. For FlexCode, KCDE and NNCDE, we solve $\hat{F}_{Y|x}(q) = \tau$ for q to obtain the estimated τ th conditional quantile. We use a test data set $\{x_1, \dots, x_k\}$ of size $k = 2000$. We consider the MSE of the estimated conditional quantile of level τ defined as

$$\text{MSE}(\tau) = \frac{1}{k} \sum_{i=1}^k [\hat{F}_{Y|X}^{-1}(\tau|X = x_i) - F_{Y|X}^{-1}(\tau|X = x_i)]^2.$$

We consider five levels of $\tau = 0.05, 0.25, 0.5, 0.75, 0.95$.

The simulation is repeated 10 times. The average MSEs and the corresponding simulation standard errors are included in Table A.1. We see all methods have larger MSEs in the tail area than in the center. This problem is more severe for KCDE and NNCDE. In model (M4), FlexCode fits the data near the first mode by using more basis functions, which leads to over-fitting in right tail area. Thus the MSE of FlexCode is small when $\tau = 0.05, 0.25$ but large when $\tau = 0.75, 0.95$. In model (M4), GCDS has a larger MSE at $\tau = 0.5$. The reason is that the estimate of GCDS is a bimodal conditional distribution with few observations around zero, which can be seen from Figure 1 that the true conditional distribution has low density around zero.

	GCDS	NNKCDE	CKDE	FlexCode
Model (M1)				
$\tau = 0.05$	0.356 (0.026)	1.875(0.011)	1.231(0.015)	1.327(0.010)
$\tau = 0.25$	0.281 (0.018)	1.355(0.005)	0.646(0.011)	0.996(0.012)
$\tau = 0.50$	0.263 (0.015)	1.028(0.009)	0.386(0.009)	0.706(0.011)
$\tau = 0.75$	0.268 (0.014)	0.800(0.013)	0.306(0.008)	0.469(0.008)
$\tau = 0.95$	0.306 (0.015)	1.106(0.018)	0.698(0.016)	0.563(0.018)
Model (M2)				
$\tau = 0.05$	1.044 (0.057)	5.318(0.028)	3.797(0.047)	4.149(0.143)
$\tau = 0.25$	0.427 (0.024)	5.002(0.023)	2.315(0.051)	3.405(0.052)
$\tau = 0.50$	0.334 (0.019)	4.818(0.035)	1.904(0.062)	3.086(0.049)
$\tau = 0.75$	0.411 (0.025)	4.600(0.046)	2.034(0.065)	2.635(0.097)
$\tau = 0.95$	0.993 (0.037)	5.395(0.104)	4.391(0.077)	3.523(0.149)
Model (M3)				
$\tau = 0.05$	0.994(0.091)	1.226(0.083)	28.51(0.729)	0.503 (0.096)
$\tau = 0.25$	1.209(0.117)	0.525(0.026)	32.50(1.371)	0.270 (0.005)
$\tau = 0.50$	8.298(0.631)	4.125 (0.134)	57.24(1.578)	4.551(0.164)
$\tau = 0.75$	9.364 (0.806)	13.52(0.205)	78.60(2.065)	19.05(0.708)
$\tau = 0.95$	17.31 (0.719)	58.25(1.495)	153.3(3.865)	73.26(2.877)
Model (M4)				
$\tau = 0.05$	0.080 (0.016)	0.319(0.004)	0.235(0.009)	0.082(0.008)
$\tau = 0.25$	0.037 (0.006)	0.150(0.002)	0.165(0.009)	0.063(0.005)
$\tau = 0.50$	0.881(0.048)	0.314(0.002)	0.466(0.010)	0.406 (0.010)
$\tau = 0.75$	0.026 (0.005)	0.156(0.002)	0.171(0.008)	0.070(0.002)
$\tau = 0.95$	0.072 (0.018)	0.313(0.004)	0.217(0.004)	0.086(0.007)

Table A.1: Mean squared prediction error (MSE) of conditional quantiles in models (M1) to (M4) and the corresponding simulation standard error(in parentheses). The smallest MSEs are in bold font.

A.1.2 The two-dimensional helix model

We now visualize some simulation results of GCDS and compare with NNKCDE and CKDE for a two-dimensional Y . FlexCode does not support the case when the dimension of Y is greater than 1, so we do not include it here.

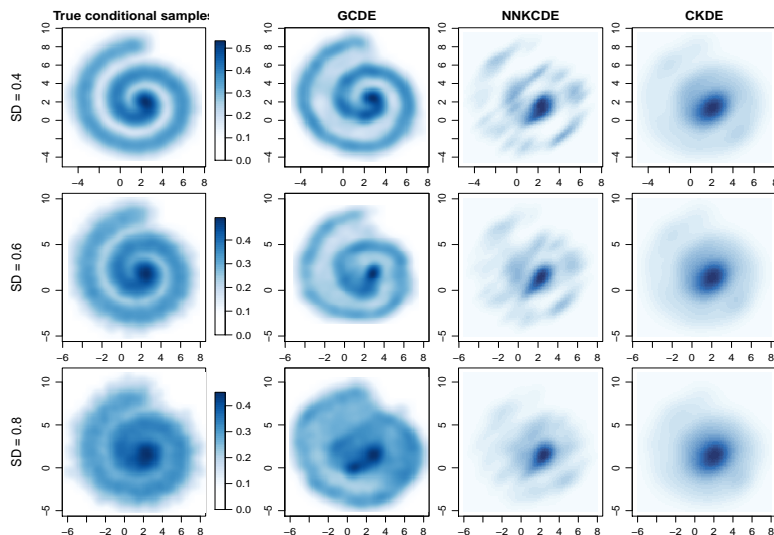


Figure A.1: Comparison of conditional density estimation in the 2D conditional helix model. From top to bottom, each row exhibits the conditional density estimates at $X = 1$ with the noise standard deviation $\sigma = 0.4, 0.6$, and 0.8 , respectively.

We consider the conditional helix model:

$$Y_1 = 2X + U \sin(2U) + \varepsilon_1,$$

$$Y_2 = 2X + U \cos(2U) + \varepsilon_2,$$

where X, U, ε_1 and ε_2 are independent, $X \sim N(0, 1), U \sim \text{Uniform}[0, 2\pi], \varepsilon_1 \sim N(0, \sigma^2)$ and $\varepsilon_2 \sim N(0, \sigma^2)$. In this model, $(U, \varepsilon_1, \varepsilon_2)$ contributes to the noise terms and is not observed.

The value of σ determines the noise level. The conditional distribution is a helix with random noise.

We visualize the quality of conditional samples and conditional density estimation given $X = 1$ at three noise levels, $\sigma = 0.4, 0.6,$ and 0.8 . The results are shown in Figure A.1. We see clearly that GCDS yields the best estimation of the conditional densities.

A.2 Effects of the neural network structure

A.2.1 Simulated data

We consider three neural network models for the generator G and another three corresponding neural network models for the discriminator D . Each set of the three networks include *Half*, *Base* and *Double*, where the size of *Half* network is half of the *Base*, and the size of *Double* is about double of *Base*. The detailed specifications of these networks are given in Table A.2.1. The *Base* network is used in the simulation studies reported in the manuscript. The number of replications in each setting is 10. The values of m are given as described above, that is, for the simulation models (M1) to (M3), we set $m = 3$; for (M4), $m = 4$.

Model	<i>Half</i>		<i>Base</i>		<i>Double</i>	
	G	D	G	D	G	D
(M1)-(M3)	$L = 1$ $W = (25)$	$L = 2$ $W = \{25, 13\}$	$L = 1$ $W = (50)$	$L = 2$ $W = (50, 25)$	$L = 2$ $W = (100, 50)$	$L = 3$ $W = (100, 50, 25)$
(M4)	$L = 2$ $W = (20, 8)$	$L = 2$ $W = \{25, 13\}$	$L = 2$ $W = (40, 15)$	$L = 2$ $W = (50, 25)$	$L = 3$ $W = (80, 30, 30)$	$L = 3$ $W = (100, 50, 25)$

Table A.2: Simulated data from models (M1)-(M4): network parameter specifications. L : the number of hidden layers, W : the widths of the layers.

Table A.3 shows the results. Since model (M3) has a larger range of distribution and heavy tail, it is more sensitive to the choice of width and depth. However, this may be

relieved by some simple transformation(For example, take logarithm). For other models, there will be some increase in MSE if the width and depth is not optimal, but the difference is mild.

	<i>Half</i>	<i>Base</i>	<i>Double</i>
(M1)	0.187(0.007)	0.152(0.005)	0.240(0.014)
(M2)	0.410(0.025)	0.333(0.024)	0.717(0.058)
(M3)	11.583(8.515)	4.031(0.425)	120.046(7.585)
(M4)	0.018(0.005)	0.019(0.003)	0.042(0.008)

Table A.3: MSE for estimating the conditional mean in the models (M1)-(M4) by GCDS, using the three neural network models *Half*, *Base* and *Double*.

A.2.2 MNIST data: effects of network structures and noise dimension m

To study the influence of the network parameters (width and depth) on the quality of the generated images with the MNIST dataset, we compare three models: *Small*, *Median* and *Large*. The network parameters (depth, width) are given in Table A.2.2.

<i>Small</i>		<i>Median</i>		<i>Large</i>	
<i>G</i>	<i>D</i>	<i>G</i>	<i>D</i>	<i>G</i>	<i>D</i>
$L = 2$	$L = 2$	$L = 3$	$L = 3$	$L = 4$	$L = 4$
$W = (128, 256)$	$W = (256, 128, 1)$	$W = (128, 256, 256)$	$W = (256, 128, 1)$	$W = (128, 128, 256, 256)$	$W = (256, 128, 128, 1)$

Table A.4: MNIST dataset: network parameter specifications. L : the number of hidden layers, W : the widths of the layers.

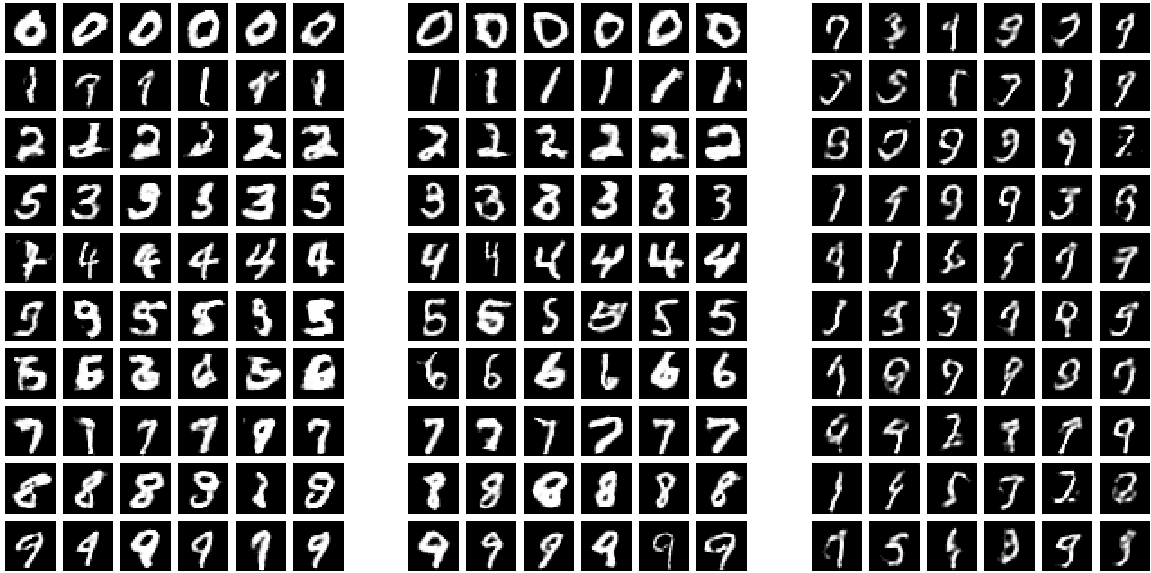


Figure A.2: Generated images using the *Half*, *Base* and *Double* neural networks. Left panel: *Half*, middle panel: *Base*, right panel: *Double*.

Figure A.2 shows three panels of images generated by GCDS using *Small*, *Median* and *Large* models from left to right. A small network has low capacity to capture the image distribution so the quality of the images are poor. Too large a network has convergence issues so the right panel does not show meaningful images.

We also examine the effects of three different values of the noise dimension m . We consider $m = 10, 100, 200$ in GCDS for generating digital images with the MNIST dataset.

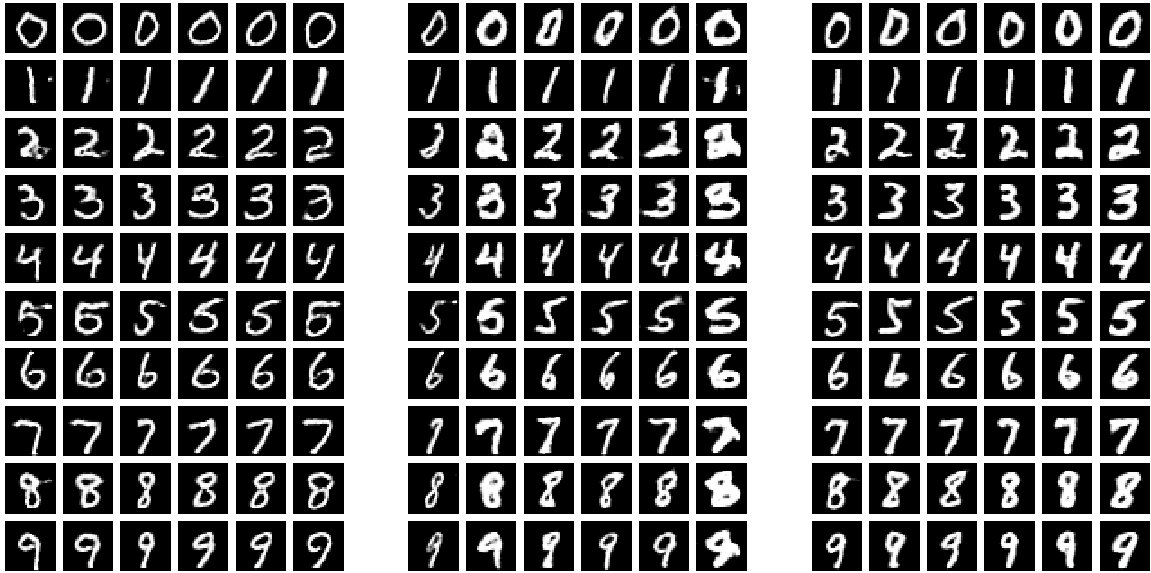


Figure A.3: Generated images from GCDS with $m = 10$ (left panel), $m = 100$ (middle panel), and $m = 200$ (right panel)

Figure A.3 shows the generated images from GCDS trained with three different values of m . From left to right, the panels correspond to the $m = 10, 100$ and 200 , respectively. The quality of the all the images look reasonable. However, when $m = 10$, there is very little variation among the images for a given digit. This means that the learned conditional distribution of the image given the label is essentially degenerate. The generated images with $m = 100$ and 200 are similar. They also show reasonable variations in the generated images. This example suggests that for the complex task of image generation with high-dimensional image data, the value of m should not be too small, and as long as m takes a reasonable large value (e.g., no smaller than the intrinsic data dimension), the results are stable. Of course, for different types of image data, the value of m should be determined on a case-by-case basis.

A.3 Performance of GCDS as sample size changes

We conduct simulation studies to demonstrate that the performance of GCDS improves as the sample increases, which supports the consistency result we obtained and also indicate reasonable convergence properties of GCDS.

In the first set of experiments, we consider models (M1)-(M4). We repeat our method with different sample sizes $n = 1000, 2500, 5000$ and 7500 . For every model, we run simulations 10 times. We record how the mean squared errors of the estimated conditional mean $E(Y|X)$ changes with sample size. The results are shown in Table A.5. There is a clear trend that the MSEs decrease with n .

n	1000	2500	5000	7500
M1	0.432(0.023)	0.262(0.012)	0.240(0.015)	0.234(0.013)
M2	2.148(0.204)	0.549(0.058)	0.308(0.024)	0.281(0.015)
M3	126.082(10.394)	23.496(12.155)	3.219(0.207)	2.339(0.136)
M4	0.060(0.008)	0.036(0.007)	0.018(0.002)	0.016(0.003)

Table A.5: MSE of the estimated condition expectation using GCDS trained with different sample sizes (simulation standard errors in parentheses).

In the second set of experiments, we examine how the samples size affects the generated images with the MNIST dataset. In this dataset, the total training set size is 60,000. We train GCDS with sample sizes $n = 5,000, 10,000$ and $60,000$.

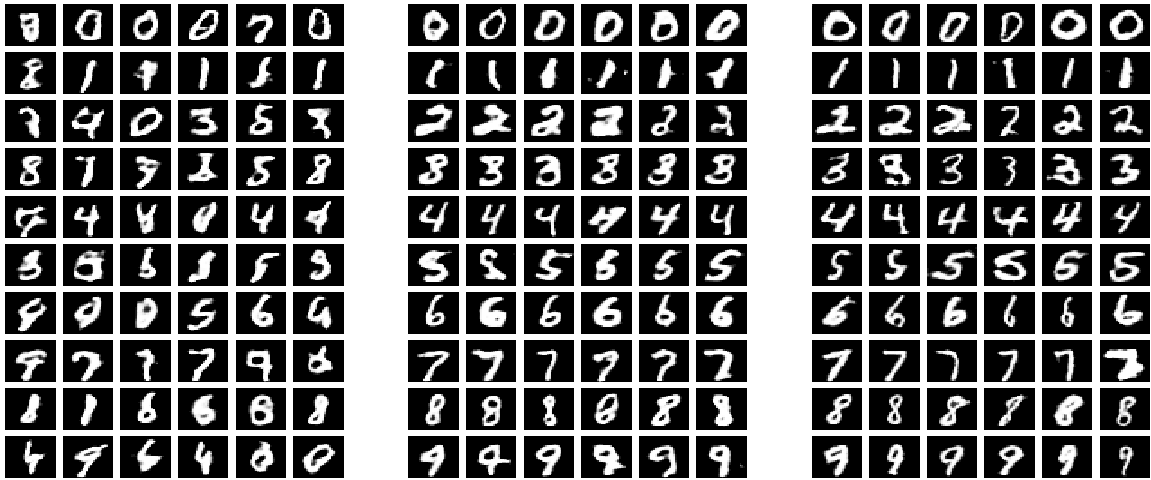


Figure A.4: From the left panel to the right panel, generated images using GCDS trained with sample sizes $n = 5,000, 10,000,$ and $60,000$.

Figure A.4 shows the generated images. From left to right, the panels represent the generated images using GCDS trained with $n = 5,000, 10,000$ and $60,000$ images in the MNIST dataset. The rows consist of generated images conditioning on the labels from 0 to 9. In the left panel consisting of images with sample size 5,000, we see that the quality of the images is poor. The quality of the images improves as the sample size increases to 10,000 and then 60,000.

A.4 Image reconstruction with STL-10 dataset

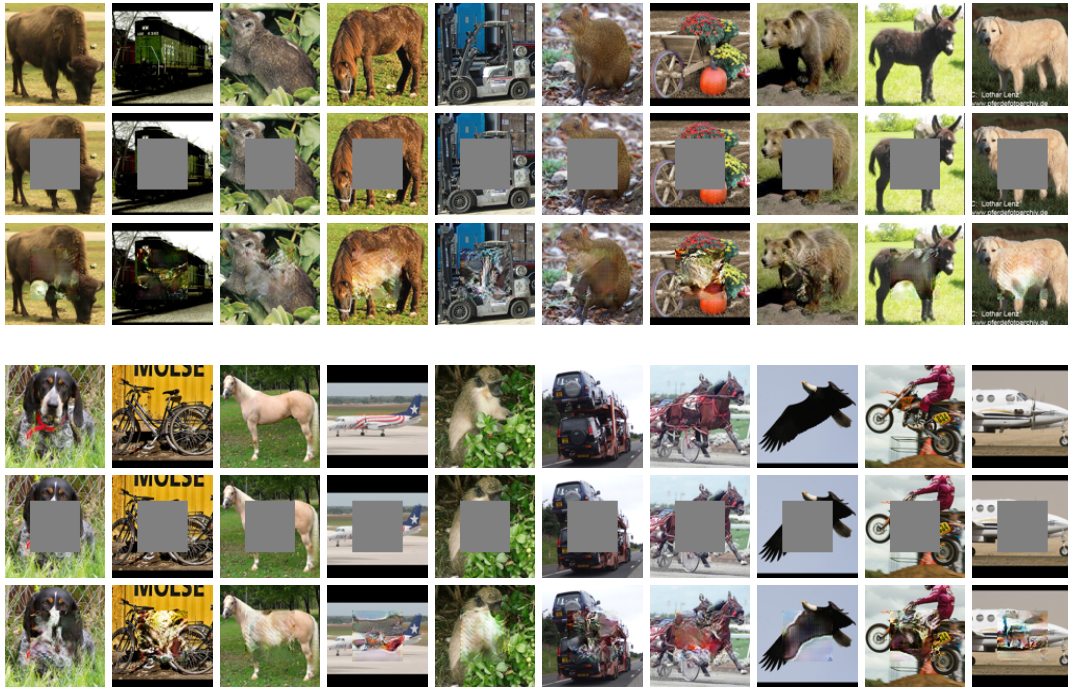


Figure A.5: Reconstructed images given partial image in STL10 dataset. In each panel, the first row consists of the true images, the second row consists of images with the corrupted area. The third row contains the reconstructed images.

As a proof of concept, we illustrate GCDS with a different high-dimensional data problem. We use the STL-10 (Coates et al., 2011) dataset, which contains color images of size $96 \times 96 \times 3$. We apply GCDS to the image reconstruction problem with the STL-10 dataset, where the central part of the image is deliberately corrupted. The size of the corrupted area is a quarter of the whole image. The task is to reconstruct the corrupted area. This is also known as image inpainting. We formulate this task as a problem of generating from a high-dimensional conditional distribution. With this formulation, the observed part of the

image is the given information while we would like to generate the missing area given the observed part of the image. For the STL-10 dataset, the response Y is the corrupted part of the image, whose dimension is $48 \times 48 \times 3 = 6,912$. The predictor X is the observed part of the image, whose dimension is $96 \times 96 \times 3 - 48 \times 48 \times 3 = 20,736$. We adopted a neural network with a structure similar to that of Iizuka et al. (2017). In each panel of Figure A.5, the first row consists of the true images, the second row shows the images with corrupted area, and the third row contains the reconstructed images. We can see that the reconstructed images are of good quality comparing with the original images. This suggests that GCDS is a promising framework for image reconstruction tasks. However, more work is needed to apply GCDS to more complex image reconstruction problems, for example, if there are multiple corrupted areas in an image. This is beyond the scope of the present work and will be studied in the future.

B Proofs and additional technical details

In this section, we give the proofs of the results stated in the paper and present additional details about f -divergence.

B.1 Proof of Lemma 3.1

Proof. Our proof follows Keziou (2003). Since $f(t)$ is convex, then for $t \in \mathbb{R}$, 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 obtain (9). This completes the proof. \square

B.2 Proof of Lemma 2.1

Proof. First, based on the basic noise-outsourcing lemma (Theorem 5.10 in Kallenberg (2002), Lemma 3.1 in Austin (2015)), there is a uniform random variable $u \sim \text{Uniform}[0, 1]$ and a measurable function $G_1 : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$ such that

$$Y = G_1(u, X) \text{ almost surely.}$$

For a random vector $\eta \sim N(\mathbf{0}, \mathbf{I}_m)$ for some $m \geq 1$, there exists a measurable function $G_2 : \mathbb{R}^m \rightarrow [0, 1]$ such that

$$u = G_2(\eta) \text{ almost surely.}$$

For example, we can simply take $G_2(\eta) = \Phi^{-1}(\eta_1)$, where Φ^{-1} is the inverse of the standard normal distribution function and η_1 is the first component of η . It is well-known that $\Phi^{-1}(\eta_1) \sim \text{Uniform}[0, 1]$. Combining the above two equations, we have

$$Y = G_1(G_2(\eta), X) \text{ almost surely.}$$

So we can simply take G as

$$G(\eta, x) = G_1(G_2(\eta), x), \quad (\eta, x) \in \mathbb{R}^m \times \mathcal{X}.$$

Therefore, we have

$$(X, G(\eta, X)) = (X, Y) \text{ almost surely.}$$

This completes the proof. □

B.3 f -divergence, convex dual and variational form

First, we note that the definition of f -divergence used in this work differs from that in Nguyen et al. (2010). The definition of f -divergence we used as defined in (7) is

$$\mathbb{D}_f(q||p) = \int p(z) f\left(\frac{q(z)}{p(z)}\right) dz. \quad (\text{C.1})$$

The definition in Nguyen et al. (2010) is

$$\mathbb{D}_\phi(q||p) = \int q(z) \phi\left(\frac{p(z)}{q(z)}\right) dz, \quad (\text{C.2})$$

where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is a convex and lower-semicontinuous function. Since

$$\mathbb{D}_f(q||p) = \int p(z) f\left(\frac{q(z)}{p(z)}\right) dz = \int q(z) \frac{p(z)}{q(z)} f\left(\frac{q(z)}{p(z)}\right) dz,$$

the two definitions (C.1) and (C.2) are equivalent if we set $f(x) = x\phi(1/x)$. For example, for the KL divergence, in definition (C.1), we take $f(x) = x \log x$, then

$$\mathbb{D}_f(q\|p) = \int p(z) \frac{q(z)}{p(z)} \log \left(\frac{q(z)}{p(z)} \right) dz = \int q(z) \log \left(\frac{q(z)}{p(z)} \right) dz = \mathbb{D}_{\text{KL}}(q\|p).$$

In definition (C.2), we take $\phi(x) = \log(1/x)$, which yields the same expression as above, i.e.,

$$\mathbb{D}_\phi(q\|p) = \int q(z) \log \left(\frac{q(z)}{p(z)} \right) dz = \mathbb{D}_{\text{KL}}(q\|p).$$

For these two functions $f(x) = x \log x$ and $\phi(x) = \log(1/x)$, we have $f(x) = x\phi(1/x)$.

We now give the expressions of the f -functions, their dual forms f^* and the variational forms of the KL-, JS- and χ^2 -divergences. For detailed derivations, we refer to Rockafellar (1970), Keziou (2003), and Nguyen et al. (2010).

- KL-divergence:

$$\begin{aligned} f(x) &= x \log x, \quad f^*(t) = e^{t-1}, \\ D_{\text{KL}}(q\|p) &= \sup_D \{ \mathbb{E}_{Z \sim q} D(Z) - \mathbb{E}_{W \sim p} e^{D(W)-1} \} \\ &= \sup_D \{ \mathbb{E}_{Z \sim q} D(Z) - \mathbb{E}_{W \sim p} e^{D(W)} \} + 1. \end{aligned}$$

- JS-divergence:

$$\begin{aligned} f(x) &= -(x+1) \log \frac{1+x}{2} + x \log x, \quad f^*(t) = -\log(2 - \exp(t)), \\ D_{\text{JS}}(q\|p) &= \sup_D \{ \mathbb{E}_{Z \sim q} D(Z) + \mathbb{E}_{W \sim p} \log(2 - \exp(D(W))) \} \\ &= \sup_D \{ \mathbb{E}_{Z \sim q} \log D(Z) + \mathbb{E}_{W \sim p} \log(1 - D(W)) \} + \log 4. \end{aligned}$$

where the last equality is obtained by change of variable: $D \rightarrow \log(2D)$. This is the form (without the constant $\log 4$) used in the GAN objective function (Goodfellow et al., 2014).

- χ^2 -divergence:

$$\begin{aligned} f(x) &= (x - 1)^2, \quad f^*(t) = t + \frac{t^2}{4}, \\ D_{\chi^2}(q||p) &= \sup_D \{ \mathbb{E}_{Z \sim q} D(Z) - \mathbb{E}_{W \sim p} [D(W) + \frac{D^2(W)}{4}] \} \\ &= \sup_D \{ 2\mathbb{E}_{Z \sim q} D(Z) - \mathbb{E}_{W \sim p} D^2(W) \} - 1, \end{aligned}$$

where the last equality is obtained by change of variable $D \rightarrow 2(D - 1)$.

B.4 Proof of Theorem 4.1

For ease of reference, we first restate the notation and Theorem 4.1 before giving the proof.

Recall $\mathcal{L}(G, D)$ is defined as

$$\mathcal{L}(G, D) = \mathbb{E}_{(X, \eta) \sim P_X P_\eta} D(X, G(\eta, X)) - \mathbb{E}_{(X, Y) \sim P_{X, Y}} \exp(D(X, Y)). \quad (\text{C.3})$$

For any measurable function $G : \mathbb{R}^m \times \mathbb{R}^d \mapsto \mathbb{R}^q$, define

$$\mathbb{L}(G) = \sup_D \mathcal{L}(G, D). \quad (\text{C.4})$$

Let

$$\widehat{\mathcal{L}}(G, D) = \frac{1}{n} \sum_{i=1}^n D(X_i, G(\eta_i, X_i)) - \frac{1}{n} \sum_{i=1}^n \exp(D(X_i, Y_i)). \quad (\text{C.5})$$

For a fixed G , let p_{XG} be the joint density of $(X, G(\eta, X))$. Lemma 3.1 implies that the optimal D is

$$D^*(z) = \log \frac{p_{XG}(z)}{p_{XY}(z)} = \log r(z).$$

Thus the optimal discriminator is the log-likelihood ratio serving as a critic of the resemblance between p_{XY} and p_{XG} . Substituting this expression into (C.4), we have,

$$\mathbb{L}(G) = \mathbb{E}_{(X,\eta) \sim P_X P_\eta} [\log r(X, G(\eta, X))].$$

Therefore, the optimal G^* minimizing the KL-divergence satisfies $P_{(X,G^*(X,\eta))} = P_{X,Y}$ by Lemma 2.2.

We make the following assumptions.

(A1) The target conditional generator $G^* : \mathbb{R}^m \times \mathcal{X} \rightarrow \mathcal{Y}$ is continuous with $\|G^*\|_\infty \leq C_0$ for some constant $0 < C_0 < \infty$.

(A2) For any $G \in \mathcal{G} \equiv \mathcal{G}_{\mathcal{D}, \mathcal{W}, \mathcal{S}, \mathcal{B}}$, $r_G(z) = p_{X,G(\eta,X)}(z)/p_{X,Y}(z) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is continuous and $0 < C_1 \leq r_G(z) \leq C_2$ for some constants $0 < C_1 \leq C_2 < \infty$.

We also make the following assumptions on the network parameters of the conditional generator G_θ and the discriminator D_ϕ .

(N1) The network parameters of \mathcal{G} satisfies

$$\mathcal{HW} \rightarrow \infty \quad \text{and} \quad \frac{\mathcal{BSH} \log(\mathcal{S}) \log n}{n} \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

(N2) The network parameters of \mathcal{D} satisfies

$$\tilde{\mathcal{H}}\tilde{\mathcal{W}} \rightarrow \infty \quad \text{and} \quad \frac{\tilde{\mathcal{B}}\tilde{\mathcal{S}}\tilde{\mathcal{H}} \log(\tilde{\mathcal{S}}) \log n}{n} \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Theorem 4.1 *Suppose that the assumptions (A1) and (A2) hold. If the network parameters of \mathcal{G} and \mathcal{D} satisfies the specifications (N1) and (N2), then*

$$\mathbb{E}_{(X_i, Y_i, \eta_i)_{i=1}^n} \|p_{X, \hat{G}_\theta(\eta, X)} - p_{X, Y}\|_{L_1}^2 \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

For simplicity, we only prove Theorem 4.1 for $q = 1$. By a truncation argument, we only need to consider on domains $E_1 = [-B, B]^{d+m} \in \mathbb{R}^d \times \mathbb{R}^m$ and $E_2 = [-B, B]^{d+q} \in \mathbb{R}^d \times \mathbb{R}^q$ with $B = \log n$. Let $C_3 = \max\{|\log(C_1)|, |\log(C_2)|\} + 1$ and $C_4 = \max_{|s| \leq 2C_3} \exp(s) + 2C_3 = \exp(2C_3) + 2C_3$. Define the empirical loss function at the sample level by

$$\hat{\mathbb{L}}(G) = \hat{\mathbb{D}}(p_{X, G(\eta, X)}, p_{X, Y}) = \sup_{D_\phi \in \mathcal{D}} \hat{\mathcal{L}}(G, D_\phi), \quad (\text{C.6})$$

where $\hat{\mathcal{L}}(G, D)$ is defined in (15). Then,

$$\hat{G}_\theta \in \arg \min_{G_\theta \in \mathcal{G}} \hat{\mathcal{L}}(G_\theta, \hat{D}_\phi),$$

and

$$\hat{D}_\phi \in \arg \max_{D_\phi \in \mathcal{D}} \hat{\mathcal{L}}(\hat{G}_\theta, D_\phi).$$

To shorten the notation, we use \hat{G} , \hat{D} , \mathcal{G} , and \mathcal{D} to denote \hat{G}_θ , \hat{D}_ϕ , $\mathcal{G}_{\mathcal{H}, \mathcal{W}, S, B}$, and $\mathcal{D}_{\tilde{\mathcal{H}}, \tilde{\mathcal{W}}, \tilde{\mathcal{S}}, \tilde{\mathcal{B}}}$,

respectively.

We first give a high-level description of the proof. By Lemma 3.2, $\mathbb{L}(G^*) = 0$. Then, by Pinsker's inequality (Tsybakov, 2008), we have

$$\|p_{X,\hat{G}(\eta,X)} - p_{X,Y}\|_{L^1}^2 \leq 2(\mathbb{L}(\hat{G}) - \mathbb{L}(G^*)). \quad (\text{C.7})$$

Next we show that the right side in (C.7) goes to 0.

For any $\bar{G} \in \mathcal{G}$, we decompose the right side in (C.7) as follows:

$$\begin{aligned} \mathbb{L}(\hat{G}) - \mathbb{L}(G^*) &= \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) \\ &\quad + \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\hat{G}, D) \end{aligned} \quad (\text{C.8})$$

$$+ \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\bar{G}, D) \quad (\text{C.9})$$

$$+ \sup_{D \in \mathcal{D}} \widehat{\mathcal{L}}(\bar{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\bar{G}, D) \quad (\text{C.10})$$

$$+ \sup_{D \in \mathcal{D}} \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(\bar{G}, D) \quad (\text{C.11})$$

$$+ \sup_D \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(G^*, D).$$

Since the terms in (C.9) and (C.11) are nonpositive, and the terms in (C.8) and (C.10) are smaller than

$$\sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)|,$$

we have

$$\begin{aligned} \mathbb{L}(\hat{G}) - \mathbb{L}(G^*) &\leq \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) + 2 \sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)| \\ &\quad + \sup_D \mathcal{L}(\bar{G}, D) - \sup_D \mathcal{L}(G^*, D). \end{aligned}$$

Taking infimum on \bar{D} over \mathcal{D} on both side of the above display we get

$$\mathbb{L}(\hat{G}) - \mathbb{L}(G^*) \leq \Delta_1 + \Delta_2 + \Delta_3, \quad (\text{C.12})$$

where

$$\begin{aligned} \Delta_1 &= \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D), \\ \Delta_2 &= 2 \sup_{G \in \mathcal{G}, D \in \mathcal{D}} |\mathcal{L}(G, D) - \hat{L}(G, D)|, \\ \Delta_3 &= \inf_{\bar{G} \in \mathcal{G}} [\mathbb{L}(\bar{G}) - \mathbb{L}(G^*)]. \end{aligned}$$

The first and the third terms Δ_1 and Δ_3 are the approximation errors and the second term Δ_2 is the statistical error. To prove the theorem, it suffices to show that these error terms converge to zero.

We first show the following Lemmas B.1 - B.3 to bound these two types of errors in (C.12).

Lemma B.1. $\Delta_3 \equiv \inf_{\bar{G} \in \mathcal{G}} [\mathbb{L}(\bar{G}) - \mathbb{L}(G^*)] = o(1)$, as $n \rightarrow \infty$.

Proof. By the assumption, G^* is continuous on $E_1 = [-B, B]^{d+m}$ with $B = \log n$, and $\|G^*\|_{L^\infty} \leq C_0$. Setting $L = \log n$, $N = n^{\frac{d+m}{2(2+d+m)}} / \log n$, $E = E_1$ and $R = B$, in Lemma B.5, we get an ReLU network $\bar{G}_{\bar{\theta}} \in \mathcal{G}$ with depth $\mathcal{H} = 12 \log n + 14 + 2(d+m)$, width $\mathcal{W} = 3^{d+m+3} \max\{(d+m)(n^{\frac{d+m}{2(2+d+m)}} / \log n)^{\frac{1}{d+m}}, n^{\frac{d+m}{2(2+d+m)}} / \log n + 1\}$, and size $\mathcal{S} = n^{\frac{d+m-2}{d+m+2}} / (\log^4 n)$, $\mathcal{B} = 2C_0$ such that

$$\|G^* - \bar{G}_{\bar{\theta}}\|_{L^\infty(E_1)} \leq 19\sqrt{d+m}\omega_f^{E_1} \left(2(\log n)n^{\frac{-1}{2+d+m}} \right). \quad (\text{C.13})$$

Let $\bar{D} = \log \frac{p_{X, \bar{G}}(\eta, X)(z)}{p_{X, Y}(z)}$ and $D^* = \log \frac{p_{X, G^*}(\eta, X)(z)}{p_{X, Y}(z)}$. Then, the above display on $\|G^* - \bar{G}\|$ in (C.13) and the continuity implies that $\|D^* - \bar{D}\| \rightarrow 0$ as $n \rightarrow \infty$. Therefore, by Lemma 3.1, we have

$$\mathbb{L}(\bar{G}) = \sup_D \mathcal{L}(\bar{G}, D) = \mathbb{E}_{(X, \eta) \sim P_X P_\eta} \bar{D}(X, \bar{G}(\eta, X)) - \mathbb{E}_{(X, Y) \sim P_{X, Y}} \exp(\bar{D}(X, Y))$$

converge to

$$\mathbb{L}(G^*) = \sup_D \mathcal{L}(G^*, D) = \mathbb{E}_{(X, \eta) \sim P_X P_\eta} D^*(X, G^*(\eta, X)) - \mathbb{E}_{(X, Y) \sim P_{X, Y}} \exp(D^*(X, Y))$$

as $n \rightarrow \infty$. □

Lemma B.2.

$$\Delta_2 \equiv \sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)| \leq \mathcal{O}(n^{-\frac{2}{2+d+m}} + n^{-\frac{2}{2+d+q}}). \quad (\text{C.14})$$

Proof. For any $G \in \mathcal{G}$, by assumption, $D_G(z) = \log r_G(z) = \log\left(\frac{p_{X, G}(\eta, X)(z)}{p_{X, Y}(z)}\right)$ is continuous on $E_2 = [-B, B]^{d+q}$ with $B = \log n$ and $\|D_G\|_{L^\infty} \leq C_3$. Setting $L = \log n$, $N = n^{\frac{d+q}{2(2+d+q)}} / \log n$, $E = E_2$ and $R = B$, in Lemma B.5, we get an ReLU network $\bar{D}_{\bar{\phi}} \in \mathcal{D}$ with depth $\tilde{\mathcal{H}} = 12 \log n + 14 + 2(d+q)$, width $\tilde{W} = 3^{d+q+3} \max\{(d+q)(n^{\frac{d+q}{2(2+d+q)}} / \log n)^{\frac{1}{d+q}}, n^{\frac{d+q}{2(2+d+q)}} / \log n + 1\}$, and size $\tilde{\mathcal{S}} = n^{\frac{d+q-2}{d+q+2}} / (\log^4 n)$, $\tilde{\mathcal{B}} = 2C_3$ such that

$$\|D_G - \bar{D}_{\bar{\phi}}\|_{L^\infty(E_2)} \leq 19\sqrt{d+m}\omega_f^{E_2} (2(\log n)n^{\frac{-1}{2+d+q}}). \quad (\text{C.15})$$

Let $Z = (X, Y) \sim P_{X, Y}$ and $Z_i = (X_i, Y_i), i = 1, \dots, n$ are i.i.d copies of Z . Let $\eta \sim P_\eta$ and $\eta \perp\!\!\!\perp X, \eta_j, j = 1, \dots, n$ are i.i.d copies of η . Then $W_i = (X_i, \eta_i)$ are i.i.d copies of $W = (X, \eta) \sim$

$P_X P_\eta$. Let $S = (W, Z) \sim (P_X P_\eta) \otimes P_{X,Y}$ and let $S_i = (W_i, Z_i) = ((X_i, \eta_i), (X_i, Y_i)), i = 1, \dots, n$ be n i.i.d copies of S . Denote

$$b(G, D; S) = D(X, G(\eta, X)) - \exp(D(X, Y)).$$

Then

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

and

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

Let $\epsilon_i, i = 1, \dots, n$ be i.i.d Rademacher random samples that are independent of $S_i, i = 1, \dots, n$. Denote the Rademacher complexity of $\mathcal{D} \times \mathcal{G}$ (Bartlett and Mendelson, 2002) by

$$\mathcal{C}(\mathcal{D} \times \mathcal{G}) = \frac{1}{n} \mathbb{E}_{\{\epsilon_i, S_i\}_{i=1}^n} \left[\sup_{G \in \mathcal{G}, D \in \mathcal{D}} \left| \sum_{i=1}^n \epsilon_i b(G, D; S_i) \right| \right].$$

Let $\mathfrak{C}(\mathcal{D} \times \mathcal{G}, e_{n,1}, \delta)$ be the covering number of $\mathcal{D} \times \mathcal{G}$ with respect to the empirical distance

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

First, by the standard symmetrization technique and the law of iterated expectations, we have

$$\begin{aligned} & \sup_{D \in \mathcal{D}, G \in \mathcal{G}} |\mathcal{L}(G, D) - \widehat{\mathcal{L}}(G, D)| \\ &= 2\mathcal{C}(\mathcal{D} \times \mathcal{G}) \\ &= 2\mathbb{E}_{S_1, \dots, S_n} \left\{ \mathbb{E}_{\epsilon_i, i=1, \dots, n} [\mathcal{C}(\mathcal{G} \times \mathcal{D}) | (S_1, \dots, S_n)] \right\}. \end{aligned} \tag{C.16}$$

For $\delta > 0$, let $\mathcal{D}_\delta \times \mathcal{G}_\delta$ be such a covering set at scale δ of $\mathcal{D} \times \mathcal{G}$. Then, by the triangle inequality and Lemma B.4 below, we have

$$\begin{aligned}
& \mathbb{E}_{S_1, \dots, S_n} \left\{ \mathbb{E}_{\epsilon_i, i=1, \dots, n} [\mathcal{C}(\mathcal{G} \times \mathcal{D}) | (S_1, \dots, S_n)] \right\} \\
& \leq \delta + \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} \left\{ \mathbb{E}_{\epsilon_i, i=1, \dots, n} \left[\sup_{(G, D) \in \mathcal{D}_\delta \times \mathcal{G}_\delta} \left| \sum_{i=1}^n \epsilon_i b(G, D; S_i) \right| \middle| (S_1, \dots, S_n) \right] \right\} \\
& \leq 2\delta + C_5 \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} \left\{ \left[\log \mathfrak{C}(\mathcal{D} \times \mathcal{G}, e_{n,1}, \delta) \right]^{1/2} \max_{(G, D) \in \mathcal{D}_\delta \times \mathcal{G}_\delta} \left[\sum_{i=1}^n b^2(G, D; S_i) \right]^{1/2} \right\}.
\end{aligned} \tag{C.17}$$

Since $\|b(G, D; S)\|_\infty \leq C_4$, we have

$$\left[\sum_{i=1}^n b^2(G, D; S_i) \right]^{1/2} \leq \sqrt{n} C_4.$$

Therefore,

$$\begin{aligned}
& \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} \left\{ (\log \mathfrak{C}(\mathcal{D} \times \mathcal{G}, e_{n,1}, \delta))^{1/2} \max_{(G, D) \in \mathcal{D}_\delta \times \mathcal{G}_\delta} \left[\sum_{i=1}^n b^2(G, D; S_i) \right]^{1/2} \right\} \\
& \leq \frac{1}{n} \mathbb{E}_{S_1, \dots, S_n} \left[(\log \mathfrak{C}(\mathcal{D} \times \mathcal{G}, e_{n,1}, \delta))^{1/2} \sqrt{n} C_4 \right] \\
& \leq \frac{C_4}{\sqrt{n}} \left[\log \mathfrak{C}(\mathcal{D}, e_{n,1}, \delta) + \log \mathfrak{C}(\mathcal{G}, e_{n,1}, \delta) \right]^{1/2}.
\end{aligned} \tag{C.18}$$

Now since $\mathfrak{C}(\mathcal{G}, e_{n,1}, \delta) \leq \mathfrak{C}(\mathcal{G}, e_{n,\infty}, \delta)$ (similar result for \mathcal{D}) and

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

where $\text{Pdim}_{\mathcal{G}}$ is the Pseudo dimension of $\mathcal{G}_{\mathcal{H},\mathcal{W},\mathcal{S},\mathcal{B}}$, which satisfies (Bartlett et al., 2019)

$$C_6 \mathcal{H} \mathcal{S} \log \mathcal{S} \leq \text{Pdim}_{\mathcal{G}} \leq C_7 \mathcal{H} \mathcal{S} \log \mathcal{S}.$$

Then, we have

$$\begin{aligned} & \frac{1}{\sqrt{n}} \left[\log \mathfrak{C}(\mathcal{D}, e_{n,1}, \delta) + \log \mathfrak{C}(\mathcal{G}, e_{n,1}, \delta) \right]^{1/2} \\ & \leq \frac{1}{\sqrt{n}} \left[\mathcal{H} \mathcal{S} \log \mathcal{S} \log \frac{\mathcal{B}n}{\delta \mathcal{H} \mathcal{S} \log \mathcal{S}} + \tilde{\mathcal{H}} \tilde{\mathcal{S}} \log \tilde{\mathcal{S}} \log \frac{\tilde{\mathcal{B}}n}{\delta \tilde{\mathcal{H}} \tilde{\mathcal{S}} \log \tilde{\mathcal{S}}} \right]^{1/2}. \end{aligned} \quad (\text{C.19})$$

Then (C.14) follows from (C.16) to (C.19) with the selection of the network parameters of $\mathcal{D}_{\tilde{\mathcal{H}},\tilde{\mathcal{W}},\tilde{\mathcal{S}},\tilde{\mathcal{B}}}$, $\mathcal{G}_{\mathcal{H},\mathcal{W},\mathcal{S},\mathcal{B}}$ and with $\delta = \frac{1}{n}$. \square

Lemma B.3. $\mathbb{E}_{(X_i, Y_i, \eta_i)_{i=1}^n} [\Delta_1] \equiv \mathbb{E}_{(X_i, Y_i, \eta_i)_{i=1}^n} [\sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D)] \rightarrow 0$ as $n \rightarrow \infty$.

Proof. Conditioning on the data $(X_i, Y_i, \eta_i)_{i=1}^n$, the supremum of $\sup_D \mathcal{L}(\hat{G}, D)$ is attained at $D_{\hat{G}}(z) = \log r_{\hat{G}}(z)$ with

$$r_{\hat{G}}(z) = \frac{p_{X, \hat{G}(\eta, X)}(z)}{p_{X, Y}(z)}.$$

By assumption, $D_{\hat{G}}(z)$ is continuous on $[-B, B]^{d+q}$ and $\|D_{\hat{G}}\|_{L^\infty} \leq C_3$. Then similar to the proof of (C.15), in Lemma B.2, there exist $\hat{D}_\phi \in \mathcal{D}$ such that $\|D_{\hat{G}} - \hat{D}_\phi\|_{L^\infty(E_2)} \rightarrow 0$, as $n \rightarrow \infty$. Then

$$0 < \sup_D \mathcal{L}(\hat{G}, D) - \sup_{D \in \mathcal{D}} \mathcal{L}(\hat{G}, D) \leq \mathcal{L}(\hat{G}, D_{\hat{G}}) - \mathcal{L}(\hat{G}, \hat{D}_\phi) \rightarrow 0,$$

by continuity. \square

Proof of Theorem 4.1

Proof. The theorem follows from (C.7) , (C.12) and Lemmas B.1 - B.3. □

Proof of Corollary 4.1

Proof. Let

$$\Delta(P_{X,\hat{G}_\theta}, P_{X,Y}) = \mathbb{E}_{X \sim P_X} \left[\int_{\mathcal{Y}} |p_{\hat{G}_\theta}(y|X) - p_{Y|X}(y|X)| dy \right].$$

We have

$$\begin{aligned} \Delta(P_{X,\hat{G}_\theta}, P_{X,Y}) &= \int_{\mathcal{X} \times \mathcal{Y}} |p_{\hat{G}_\theta}(y|x) - p_{Y|X}(y|x)| p_X(x) dx dy \\ &= \int_{\mathcal{X} \times \mathcal{Y}} |p_{X,\hat{G}_\theta}(x, y) - p_{X,Y}(x, y)| dx dy. \end{aligned}$$

It follows from Theorem 4.1 that

$$\mathbb{E}_{\{X_i, Y_i\}_{i=1}^n} \Delta(P_{X,\hat{G}_\theta}(\eta, X), P_{X,Y}) \rightarrow 0.$$

So Corollary 4.1 follows from the Markov inequality. □

Finally, we state the two lemmas used in the proofs above for ease of reference, The first lemma is used in the proof of Lemma B.2.

Lemma B.4. *If $\xi_i, i = 1, \dots, m$ are m finite linear combinations of Rademacher variables $\epsilon_j, j = 1, \dots, J$. Then*

$$\mathbb{E}_{\epsilon_j, j=1, \dots, J} \max_{1 \leq i \leq m} |\xi_i| \leq C_6 (\log m)^{1/2} \max_{1 \leq i \leq m} (\mathbb{E} \xi_i^2)^{1/2} \tag{C.20}$$

for some constant $C_6 > 0$.

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

The following approximation result about neural networks (Shen et al., 2020) is used in the proof of Lemma B.1.

Lemma B.5. *Let f be a uniformly continuous function defined on $E \subseteq [-R, R]^d$. For arbitrary $L \in \mathbb{N}^+$ and $N \in \mathbb{N}^+$, there exists a function ReLU network f_ϕ with width $3^{d+3} \max\{d \lfloor N^{1/d} \rfloor, N + 1\}$ and depth $12L + 14 + 2d$ such that*

$$\|f - f_\phi\|_{L^\infty(E)} \leq 19\sqrt{d}\omega_f^E(2RN^{-2/d}L^{-2/d}),$$

where, $\omega_f^E(t)$ is the modulus of continuity of f satisfying $\omega_f^E(t) \rightarrow 0$ as $t \rightarrow 0^+$.

Proof. This is Theorem 4.3 in Shen et al. (2020). \square

References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. (2016). Tensorflow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 265–283.
- Ali, S. M. and Silvey, S. D. (1966). A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B (Methodological)*, 28(1):131–142.
- Austin, T. (2015). Exchangeable random measures. *Annales de l’I.H.P. Probabilités et statistiques*, 51(3):842–861.

- Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. (2019). Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20:1–17.
- Bartlett, P. L. and Mendelson, S. (2002). Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3:463–482.
- Bauer, B. and Kohler, M. (2019). On deep learning as a remedy for the curse of dimensionality in nonparametric regression. *Ann. Statist.*, 47(4):2261–2285.
- Bengio, Y., Courville, A., and Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828.
- Bhattacharya, P. K. and Gangopadhyay, A. K. (1990). Kernel and nearest-neighbor estimation of a conditional quantile. *Ann. Statist.*, 18(3):1400–1415.
- Bishop, C. (2006). *Pattern Recognition and Machine Learning*. Springer.
- Bott, A.-K. and Kohler, M. (2017). Nonparametric estimation of a conditional density. *Annals of the Institute of Statistical Mathematics*, 69(1):189–214.
- Chen, M., Jiang, H., Liao, W., and Zhao, T. (2019). Nonparametric regression on low-dimensional manifolds using deep relu networks. *arXiv:1908.01842*.
- Chen, X. and Linton, O. (2001). The estimation of conditional densities. In *Asymptotics in Statistics and Probability, Festschrift for George Roussas*, ed. M.L. Puri.
- Coates, A., Ng, A., and Lee, H. (2011). An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 215–223. JMLR Workshop and Conference Proceedings.
- Cook, R. (1998). *Regression Graphics: Ideas for Studying Regressions Through Graphics*. Wiley Series in Probability and Statistics. Wiley.
- Dalmasso, N., Pospisil, T., Lee, A. B., Izbicki, R., Freeman, P. E., and Malz, A. I. (2020). Conditional density estimation tools in python and r with applications to photometric redshifts and likelihood-free cosmological inference. *Astronomy and Computing*, page 100362.

- De la Pena, V. and Giné, E. (2012). *Decoupling: from Dependence to Independence*. Springer Science & Business Media.
- Dua, D. and Graff, C. (2017). UCI machine learning repository. <https://archive.ics.uci.edu/ml/index.php>.
- Fan, J., Yao, Q., and Tong, H. (1996). Estimation of conditional densities and sensitivity measures in nonlinear dynamical systems. *Biometrika*, 83(1):189–206.
- Fan, J. and Yim, T. H. (2004). A crossvalidation method for estimating conditional densities. *Biometrika*, 91(4):819–834.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in Neural Information Processing Systems*, volume 27, pages 2672–2680.
- Hall, P., Racine, J., and Li, Q. (2004). Cross-validation and the estimation of conditional probability densities. *Journal of the American Statistical Association*, 99(468):1015–1026.
- Hall, P. and Yao, Q. (2005). Approximating conditional distribution functions using dimension reduction. *Annals of Statistics*, 33(3):1404–1421.
- Hyndman, R. J., Bashtannyk, D. M., and Grunwald, G. K. (1996). Estimating and visualizing conditional densities. *Journal of Computational and Graphical Statistics*, 5(4):315–336.
- Hyvärinen, A. and Pajunen, P. (1999). Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, 12(3):429–439.
- Iizuka, S., Simo-Serra, E., and Ishikawa, H. (2017). Globally and locally consistent image completion. *ACM Transactions on Graphics (ToG)*, 36(4):1–14.
- Izbicki, R. and Lee, A. B. (2016). Nonparametric conditional density estimation in a high-dimensional regression setting. *Journal of Computational and Graphical Statistics*, 25(4):1297–1316.
- Izbicki, R., Lee, A. B., et al. (2017). Converting high-dimensional regression to high-dimensional conditional density estimation. *Electronic Journal of Statistics*, 11(2):2800–2831.

- Jiao, Y., Shen, G., Lin, Y., and Huang, J. (2021). Deep nonparametric regression on approximately low-dimensional manifolds. *arXiv 2104.06708*.
- Kallenberg, O. (2002). *Foundations of Modern Probability*. Springer-Verlag, New York, 2nd edition.
- Keziou, A. (2003). Dual representation of φ -divergences and applications. *Comptes Rendus Mathématique*, 336(10):857–862.
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representation*.
- Kohler, M. and Langer, S. (2020). On the rate of convergence of fully connected very deep neural network regression estimates. *arXiv 1908.11133*. (to appear in *Annals of Statistics*).
- LeCun, Y., Cortes, C., and Burges, C. (2010). Mnist handwritten digit database. *AT&T Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2.
- Li, K.-C. (1991). Sliced inverse regression for dimension reduction. *Journal of the American Statistical Association*, 86(414):316–327.
- Mirza, M. and Osindero, S. (2014). Conditional generative adversarial nets. cite arxiv:1411.1784.
- Nakada, R. and Imaizumi, M. (2019). Adaptive approximation and estimation of deep neural network with intrinsic dimensionality. *arXiv:1907.02177*.
- Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2010). Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861.
- Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in neural information processing systems*, pages 271–279.
- Rockafellar, T. R. (1970). *Convex analysis*. Princeton University Press.
- Rosenblatt, M. (1969). Conditional probability density and regression estimators. In *In Multivariate Analysis II, Ed. P. R. Krishnaiah*, pages 25–31, New York. Academic Press, New York.

- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117.
- Schmidt-Hieber, J. (2020). Nonparametric regression using deep neural networks with ReLU activation function (with discussion). *Ann. Statist.*, 48(4):1916–1921.
- Scott, D. W. (1992). *Multivariate Density Estimation: Theory, Practice and Visualization*. Wiley, New York.
- Shen, G., Jiao, Y., Lin, Y., Horowitz, J. L., and Huang, J. (2021a). Deep quantile regression: Mitigating the curse of dimensionality through composition. *arXiv 2107.04907*.
- Shen, G., Jiao, Y., Lin, Y., and Huang, J. (2021b). Robust nonparametric regression with deep neural networks. *arXiv 2107.10343*.
- Shen, Z., Yang, H., and Zhang, S. (2020). Deep network approximation characterized by number of neurons. *Commun. Comput. Phys.*, 28(5):1768–1811.
- Sugiyama, M., Takeuchi, I., Suzuki, T., Kanamori, T., Hachiya, H., and Okanojara, D. (2010). Least-squares conditional density estimation. *IEICE Transactions on Information and Systems*, 93(3):583–594.
- Tsybakov, A. (2008). *Introduction to Nonparametric Estimation*. Springer Science & Business Media.
- Van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Processes: with Applications to Statistics*. Springer, New York.
- Zhao, L. and Liu, Z. (1985). Strong consistency of the kernel estimators of conditional density function. *Acta Mathematica Sinica*, 1:314–318.