
An Error Analysis of Deep Density-Ratio Estimation with Bregman Divergence

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We establish non-asymptotic error bounds for a nonparametric density-ratio estimator
2 using deep neural networks with the Bregman divergence. We also show that
3 the deep density-ratio estimator can mitigate the curse of dimensionality when the
4 data is supported on an approximate low-dimensional manifold. Our error bounds
5 are optimal in the minimax sense and the pre-factors in our error bounds depend
6 on the dimensionality of the data polynomially. We apply our results to investigate
7 the convergence properties of the telescoping density-ratio estimator (Rhodes et al.,
8 2020) and provide sufficient conditions under which it has a smaller upper error
9 bound than a single-ratio estimator.

10 1 Introduction

11 Density-ratio estimation is of key importance in various statistical and machine learning problems
12 (Sugiyama et al., 2012b; Kato & Teshima, 2021). There is a vast literature on density-ratio estimation
13 due to its wide range of applications, such as discriminative analysis (Silverman, 1978; Cox & Ferry,
14 1991), covariate shift adaptation (Sugiyama et al., 2008; Tsuboi et al., 2009), two-sample testing
15 (Qin, 1998; Sugiyama et al., 2011), energy-based modelling (Gutmann & Hyvärinen, 2012; Ceylan &
16 Gutmann, 2018), generative learning (Goodfellow et al., 2014; Nowozin et al., 2016), and mutual
17 information estimation (Moustakides & Basioti, 2019; Rhodes et al., 2020), among others.

18 Let Z_q and $Z_p \in \mathcal{Z} = [0, 1]^d$ be two random vectors with probability density functions q^* and p^* ,
19 respectively. Given independent and identically distributed (i.i.d) samples $\{Z_{q,i}\}_{i=1}^{n_q}$ from q^* and
20 $\{Z_{p,j}\}_{j=1}^{n_p}$ from p^* , a basic problem is to estimate the density ratio

$$R^*(z) = q^*(z)/p^*(z), z \in \mathcal{Z}.$$

21 A naive estimator of R^* is \hat{q}/\hat{p} , where \hat{q} and \hat{p} are the density estimators of q^* and p^* , respectively.
22 However, such an estimator can be highly unstable. Moreover, density estimation itself is a difficult
23 problem, especially in the high-dimensional settings. For example, kernel density estimators
24 (Rosenblatt, 1956; Parzen, 1962) works well when $d \leq 3$, but deteriorate dramatically as d increases.
25 To avoid density estimation, various methods have been proposed to estimate the density ratio R^*
26 directly, including the density matching approach (Sugiyama et al., 2008; Tsuboi et al., 2009; Yamada
27 & Sugiyama, 2009; Nguyen et al., 2010; Yamada et al., 2010), the moment matching approach (Qin,
28 1998; Gretton et al., 2009; Kanamori et al., 2012b), the density-ratio fitting approach (Kanamori
29 et al., 2009, 2012a), and the unified density-ratio matching approach under Bregman divergence
30 framework (Sugiyama et al., 2012a). Impressive empirical successes of using deep neural networks
31 in density-ratio estimation have been reported in some recent works (Moustakides & Basioti, 2019;
32 Rhodes et al., 2020). Moreover, Kato & Teshima (2021) studied the convergence properties of deep
33 density-ratio estimation under a modified Bregman divergence criterion.

34 In this paper, we study deep density-ratio estimators with the Bregman divergence as the criterion. We
 35 apply our results to construct an estimator for statistical inference for the Kullback-Liebler divergence.
 36 We also study the theoretical properties of the telescoping density-ratio estimator (Rhodes et al.,
 37 2020) based on our results.

38 Our contributions are as follows:

- 39 1. We establish non-asymptotic error bounds for the density-ratio estimator using deep neural
 40 networks under the Bregman divergence (BD, Bregman, 1967), and provide a neural network
 41 architecture for the estimator to achieve minimax optimal rate $O_p(n^{-2\beta/(d+2\beta)})$, where
 42 $n = \min\{n_q, n_p\}$ and β is a smoothness parameter of the logarithmic density-ratio function;
 43 see Subsection 3.2 for details;
- 44 2. We show that deep density-ratio estimator with the Bregman divergence criterion is able
 45 to mitigate the curse of dimensionality when the data is supported on an approximate
 46 low-dimensional manifold; see Subsection 3.3;
- 47 3. We apply our results to study the convergence properties of the telescoping density-ratio
 48 estimator (Rhodes et al., 2020) and demonstrate its advantages over single-ratio estimators
 49 under certain conditions.

50 **Notation.** Let $n = \min\{n_q, n_p\}$ be the smaller sample size between the two samples $\{Z_{q,i}\}_{i=1}^{n_q}$
 51 and $\{Z_{p,j}\}_{j=1}^{n_p}$. In addition, $\|\cdot\|_\infty$ denotes the sup-norm on some specific domain, and C, C_0 are
 52 generic constants that may vary from place to place. For any measurable function f , we denote
 53 $\|f\|_{\max} := \max\{\|f\|_p, \|f\|_q\}$ and $\|f\|_{n_p, n_q} = \max\{\|f\|_{p, n_p}, \|f\|_{q, n_q}\}$, where $\|f\|_I^2 = E_{I^*} f^2(Z)$
 54 and $\|f\|_{I, n_I}^2 = E_{n_I} f^2(Z) = (1/n_I) \sum_{t=1}^{n_I} f^2(Z_{I,t})$, $I = p, q$.

55 2 Density-ratio estimation

56 In this section, we first present the density-ratio estimation problem using the Bregman divergence
 57 (BD, Bregman, 1967) and then describe the structure of the deep neural networks to be used in
 58 density-ratio estimation.

59 Let $\psi : \mathbb{R} \rightarrow \mathbb{R}$ be a first-order continuously differentiable and strictly convex function. Define

$$\Delta_\psi(x, y) = \psi(x) - \psi(y) - \psi'(y)(x - y),$$

60 where ψ' is the derivative of ψ . Then, the convexity of ψ implies that $\Delta_\psi(x, y) \geq 0$ and the equality
 61 holds if and only if $x = y$. It follows that $E_{p^*} \Delta_\psi(R^*(Z), R(Z)) \geq 0$ and the equality holds if and
 62 only if $R = R^*$. Therefore, the target density-ratio $R^* = q^*/p^*$ can be characterized as a minimizer:

$$R^* \in \arg \min_{R \text{ nonnegative and measurable}} E_{p^*} \Delta_\psi(R^*(Z), R(Z)).$$

63 We verify in the appendix that

$$\begin{aligned} & E_{p^*} \Delta_\psi(R^*(Z), R(Z)) \\ &= E_{p^*} [\psi'(R(Z))R(Z) - \psi(R(Z))] - E_{q^*} [\psi'(R(Z))] + E_{p^*} [\psi(R^*(Z))] \end{aligned} \quad (1)$$

64 Since the last term on the right side in (1) $E_{p^*} [\psi(R^*(Z))]$ is independent of R , we have

$$R^* \in \arg \min_{R \text{ nonnegative and measurable}} E_{p^*} [\psi'(R(Z))R(Z) - \psi(R(Z))] - E_{q^*} [\psi'(R(Z))]. \quad (2)$$

65 Hence, for any measurable function $R : \mathcal{Z} \rightarrow \mathbb{R}$, the BD score induced by ψ for estimating the target
 66 density-ratio $R^* = q^*/p^*$ is

$$\mathcal{B}_\psi(R) = E_{p^*} [\psi'(R(Z))R(Z) - \psi(R(Z))] - E_{q^*} [\psi'(R(Z))], \quad (3)$$

67 where ψ' is the derivative of ψ (Sugiyama et al., 2012a,b). Then, R^* is the minimizer of $\mathcal{B}_\psi(R)$ over
 68 all nonnegative measurable functions.

69 Because a density ratio is always nonnegative, a nonnegative constraint needs to be considered when
 70 defining the density ratio as a minimizer, as in (2). This makes the minimization problem more
 71 difficult to solve. To avoid the non-negative constraint of the density ratio, we first consider the

72 log-density ratio $D^* := \log R^*$. Then the nonnegativity constraint is no longer needed and by (2), we
 73 have

$$D^* \in \arg \min_{D \text{ measurable}} \mathcal{B}_\psi(\exp(D)).$$

74 In practice, the estimation of R^* can be based on an empirical version of \mathcal{B}_ψ when random samples
 75 from p^* and q^* are available. Suppose we have samples $\{Z_{q,i}\}_{i=1}^{n_q}$ i.i.d. q^* and $\{Z_{p,j}\}_{j=1}^{n_p}$ i.i.d. p^* .
 76 We estimate D^* by

$$\widehat{D} \in \arg \min_{D \in \mathcal{F}_n} \widehat{\mathcal{B}}_\psi(e^D), \quad (4)$$

77 where \mathcal{F}_n is a class of neural network functions and $\widehat{\mathcal{B}}_\psi(e^D)$ is an empirical version of $\mathcal{B}_\psi(e^D)$
 78 defined in (3), which can be written as

$$\widehat{\mathcal{B}}_\psi(e^D) = \frac{1}{n_p} \sum_{j=1}^{n_p} \mathcal{L}_1(D(Z_{p,j})) + \frac{1}{n_q} \sum_{i=1}^{n_q} \mathcal{L}_2(D(Z_{q,i})),$$

79 where

$$\mathcal{L}_1(t) = \psi'(e^t)e^t - \psi(e^t) \text{ and } \mathcal{L}_2(t) = -\psi'(e^t). \quad (5)$$

80 The density-ratio estimator is $\widehat{R} = \exp(\widehat{D})$.

81 We take the function class \mathcal{F}_n to be $\mathcal{F}_{M,\mathcal{D},\mathcal{W},\mathcal{U},\mathcal{S}}$, a class of ReLU activated feedforward neural
 82 networks (FNNs) $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ with parameter θ , depth \mathcal{D} , width \mathcal{W} , size \mathcal{S} , number of neurons
 83 \mathcal{U} . We require that $\|f_\theta\|_\infty \leq M$ for some $0 \leq M \leq \infty$. There are \mathcal{D} hidden layers and $(\mathcal{D} + 1)$
 84 layers in total. The width \mathcal{W} is the maximum width of the hidden layers; the number of neurons \mathcal{U} is
 85 defined as the number of neurons of f_θ ; the size \mathcal{S} is the total number of parameters in the network.
 86 Note that $\mathcal{D}, \mathcal{W}, \mathcal{U}, \mathcal{S}$ may depend on n , but we suppress the dependence for notational simplicity.
 87 We write $\mathcal{F}_{M,\mathcal{D},\mathcal{W},\mathcal{U},\mathcal{S}}$ as \mathcal{F}_{FNN} for brevity.

88 3 Theoretical results

89 In this section, we first study the error bounds for the deep logarithmic density-ratio estimator. The
 90 bounds for the density-ratio estimator follows directly based on the properties of the exponential
 91 function. We also show that deep density-ratio estimator can mitigate the curse of dimensionality
 92 when data is supported on an approximate low-dimensional manifold.

93 3.1 General error bounds

94 To state our assumptions and results, we need the definitions of μ -smoothness, σ -strong convexity
 95 and pseudo dimension.

96 A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is said to be μ -smooth over a set $\mathcal{A} \subseteq \mathbb{R}$ if it is differentiable over \mathcal{A} and its
 97 first-order derivative f' satisfies

$$|f'(x) - f'(y)| \leq \mu|x - y|, \quad \forall x, y \in \mathcal{A}, \quad (6)$$

98 where $0 \leq \mu < \infty$. The constant μ is called the smoothness parameter.

99 A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is called σ -strongly convex if the domain $\text{dom}(f)$ of f is convex and for any
 100 $x, y \in \text{dom}(f)$ and $\lambda \in [0, 1]$, f satisfies

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) - \frac{\sigma}{2}\lambda(1 - \lambda)(x - y)^2, \quad (7)$$

101 where $0 \leq \sigma < \infty$. The constant σ is called the strong convexity (SC) parameter.

102 For a function class \mathcal{F} , its pseudo dimension denoted by $\text{Pdim}(\mathcal{F})$, is the largest integer B satisfying
 103 that there exists $(x_1, x_2, \dots, x_B, y_1, y_2, \dots, y_B) \in \mathcal{Z}^B \times \mathbb{R}^B$ such that for any $(r_1, r_2, \dots, r_B) \in$
 104 $\{0, 1\}^B$, there exists an $f \in \mathcal{F}$ satisfying for any $i \in \{1, 2, \dots, B\} : f(x_i) > y_i \Leftrightarrow r_i = 1$ (Anthony
 105 & Bartlett, 1999; Bartlett et al., 2019).

Table 1: Commonly-used Loss Functions ψ

Name	$\psi(c)$	Domain	Smooth Parameter μ	SC Parameter σ
LS	$(c-1)^2$	\mathbb{R}	2	2
LR	$c \log c - (c+1) \log(c+1)$	$[a, b] (-1 \leq a \leq b)$	$\frac{1}{a(a+1)}$	$\frac{1}{b(b+1)}$
KL	$c \log c - c$	$[a, b] (0 \leq a \leq b)$	$\frac{1}{a}$	$\frac{1}{b}$

106 **Remark 1.** For any measurable function class \mathcal{F} , by the definition of VC dimension,
 107 $\text{VCdim}(\mathcal{F}) \leq \text{Pdim}(\mathcal{F})$. If \mathcal{F} is the class of functions generated by ReLU FNNs, it follows from
 108 Theorem 14.1 of Anthony & Bartlett (1999) that $\text{Pdim}(\mathcal{F}) \leq \text{VCdim}(\mathcal{F})$. Hence, for the function
 109 class \mathcal{F} generated by ReLU FNNs, $\text{Pdim}(\mathcal{F}) = \text{VCdim}(\mathcal{F})$.

110 We make the following assumptions.

111 **Assumption 1.** The function ψ is μ -smooth & σ -strongly convex, that is, it satisfies (6) and (7).

112 Some commonly-used ψ 's satisfy Assumption 1; see Table 1 for some examples.

113 **Assumption 2.** There exists a constant $0 < M < \infty$ such that $\|D^*\|_\infty \leq M, \|D\|_\infty \leq M$ for
 114 every $D \in \mathcal{F}_{\text{FNN}}$.

115 Assumption 2 assumes that the target density ratio is bounded. Such an assumption is often made in
 116 nonparametric statistics for avoiding technical difficulties associated with dealing with unbounded
 117 functions. We will partially relax this assumption below. The finite M in Assumption 2 can
 118 be relaxed to $M = \mathcal{O}(\log \log n)$ at a small price of an additional logarithm term in the error
 119 bounds. The boundedness of a network can be achieved by clipping operation. For example, let
 120 $T_M(t) = -MI\{t < -M\} + tI\{-M \leq t \leq M\} + MI\{t > M\}$ be the truncation function taking
 121 values in $[-M, M]$, then $T(t) = \sigma(t) - \sigma(\sigma(t) - M) - \{\sigma(-t) - \sigma(\sigma(-t) - M)\}$ can be computed
 122 by a ReLU network with depth 2 and width 4. Hence, through network concatenation, we can
 123 construct some bounded ReLU FNNs and such a boundedness assumption can be satisfied.

124 Define the best in class approximation of D^* in \mathcal{F}_{FNN} as $D_{\text{NN}} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \|D - D^*\|_{\max}$.
 125 Denote

$$\xi_n = \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}}. \quad (8)$$

126 **Theorem 1.** Suppose Assumptions 1-2 are satisfied. When $n \geq \text{Pdim}(\mathcal{F}_{\text{FNN}})$, there exists a
 127 constant C depending on (μ, σ, M) such that for any $\gamma > 0$, with probability at least $1 - \exp(-\gamma)$,

$$\|\widehat{D} - D^*\|_{\max} \leq C \left(\xi_n + \|D_{\text{NN}} - D^*\|_{\max} + \sqrt{\frac{\gamma}{n}} \right),$$

128 and

$$\|\widehat{D} - D^*\|_{n_p, n_q} \leq 2C \left(\xi_n + \|D_{\text{NN}} - D^*\|_{\max} + \sqrt{\frac{\gamma}{n}} \right).$$

129 We have the following corollary for the expected error.

130 **Corollary 1.** Under the conditions of Theorem 1, there exists a constant C_0 depending only on
 131 (μ, σ, M) , such that

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \leq C_0 (\xi_n^2 + \|D_{\text{NN}} - D^*\|_{\max}^2),$$

132 and

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{n_p, n_q}^2 \leq 2C_0 (\xi_n^2 + \|D_{\text{NN}} - D^*\|_{\max}^2).$$

133 The above results are obtained under the boundedness Assumption 2. While such an assumption is
 134 often made in the error analysis of nonparametric procedures, it is somewhat restrictive in density-
 135 ratio estimation problems. For example, this assumption may not be satisfied in the presence of
 136 the density-chasm problem, i.e., the gap between two densities is large (Rhodes et al., 2020). We
 137 establish an error bound result with the following partially relaxed assumption.

138 **Assumption 3.** *There exists a constant $0 < M < \infty$ such that $D^*(z) \geq -M$ for every $z \in \mathcal{Z}$*
 139 *and $\|D\|_\infty \leq M$ for every $D \in \mathcal{F}_{\text{FNN}}$.*

140 This assumption does not require the target log-density ratio D^* to be bounded above. Denote
 141 truncated versions of D^* and R^* by

$$\begin{aligned} D_M^*(z) &= D^*(z)\mathbf{1}\{D^*(z) \leq M\} + M\mathbf{1}\{D^*(z) \geq M\}, \\ R_M^*(z) &= R^*(z)\mathbf{1}\{R^*(z) \leq e^M\} + e^M\mathbf{1}\{R^*(z) \geq e^M\}, \end{aligned}$$

142 where $0 < M < \infty$ and $\mathbf{1}\{\cdot\}$ is the indicator function. We establish a non-asymptotic error bound
 143 involving the truncation error.

144 **Theorem 2.** *Suppose Assumptions 1 and 3 hold. When $n > \text{Pdim}(\mathcal{F}_{\text{FNN}})$, there exists two*
 145 *constants C depending only on (μ, σ, M) and C_0 depending only on (μ, σ) , such that*

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_p^2 \leq C_0 e^{2M} \|R^* - R_M^*\|_p^2 + C \left(\xi_n + \inf_{D \in \mathcal{F}_{\text{FNN}}} \|D - D_M^*\|_p^2 \right),$$

146 where ξ_n is defined in (8).

147 The term $\|R^* - R_M^*\|_p^2$ is the truncation error for an unbounded R^* and the unboundedness also
 148 leads to the term $\xi_n = \lceil \text{Pdim}(\mathcal{F}_{\text{FNN}}) (\log n) / n \rceil^{1/2}$ in the error bound, which is greater than ξ_n^2 in
 149 the bounded case. However, because no boundedness assumption is needed in this theorem, we can
 150 apply it to study the convergence properties of the telescoping density-ratio estimator of Rhodes et al.
 151 (2020) in Section 4 below.

152 3.2 Non-asymptotic error bounds

153 By Corollary 1, it suffices to bound the estimation error $\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n / n$ and the approxi-
 154 mation error $\|D_{\text{NN}} - D^*\|_{\max}^2$. It follows from Theorem 6 in Bartlett et al. (2019) that, for
 155 $\mathcal{F}_{\text{FNN}} = \mathcal{F}_{M, \mathcal{D}, \mathcal{W}, \mathcal{U}, \mathcal{S}}$, there exists a universal constant C_2 such that $\text{Pdim}(\mathcal{F}_{\text{FNN}}) \leq C_2 \mathcal{S} \mathcal{D} \log \mathcal{S}$.
 156 To control the approximation error $\|D_{\text{NN}} - D^*\|_{\max}^2$, we assume that D^* belongs to the Hölder
 157 class $\mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$ and $a \in (0, 1]$, where \mathbb{N}^+ is the set of positive
 158 integers.

159 **Definition 1** (Hölder class). *A Hölder class $\mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$ and*
 160 *$a \in (0, 1]$ consists of function $f : [0, 1]^d \rightarrow \mathbb{R}$ satisfying*

$$\max_{\|\alpha\|_1 \leq k} \|\partial^\alpha f\|_\infty, \max_{\|\alpha\|_1 = k} \max_{x \neq y} \frac{|\partial^\alpha f(x) - \partial^\alpha f(y)|}{\|x - y\|_2^a} \leq M,$$

161 where $\|\alpha\|_1 = \sum_{i=1}^d \alpha_i$ and $\partial^\alpha = \partial^{\alpha_1} \partial^{\alpha_2} \dots \partial^{\alpha_d}$ for $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_d) \in \mathbb{N}^{+d}$.

162 We use Theorem 3.3 of Jiao et al. (2021) to control the approximation error $\|D_{\text{NN}} - D^*\|_{\max}^2$. For
 163 convenience, we include this result in the following lemma.

164 We specify the width \mathcal{W} and depth \mathcal{D} as follows. For any $K, L \in \mathbb{N}^+$,

$$\mathcal{W} = 38(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1} L \lceil \log_2(8L) \rceil, \mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 K \lceil \log_2(8K) \rceil, \quad (9)$$

165 where $\lceil a \rceil$ is the smallest integer no less than a .

166 **Lemma 1** (Approximation error). *Assume $f \in \mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$*
 167 *and $a \in (0, 1]$. Then there exists a function ϕ_0 implemented by a ReLU network with width \mathcal{W} and*
 168 *depth \mathcal{D} specified in (9) such that*

$$\sup_{x \in [0, 1]^d \setminus H_{B, \delta}} |f - \phi_0| \leq 18MC_\beta (KL)^{-\frac{2\beta}{d}},$$

169 where $C_\beta = (\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + (\beta \vee 1)/2}$, $H_{B, \delta} = \cup_{i=1}^d \{x = [x_1, \dots, x_d] : x_i \in$
 170 $\cup_{b=1}^{B-1} (b/B - \delta, b/B)\}$ for $B = \lceil (KL)^{2/d} \rceil$, $\delta \in (0, 1/(3B))$ and $a \vee b = \max(a, b)$.

171 Furthermore, if $\mathcal{W} = 38(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1} 3^d L \lceil \log_2(8L) \rceil$, $\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 K \lceil \log_2(8K) \rceil + 2d$,
 172 then

$$\sup_{x \in [0, 1]^d} |f - \phi_0| \leq 19MC_\beta (KL)^{-\frac{2\beta}{d}}.$$

173 In this uniform approximation result, the width \mathcal{W} is required to depend on d exponentially.

174 The following theorem gives an error bound for \widehat{D} .

175 **Theorem 3** (Non-asymptotic error bound for \widehat{D}). *Suppose that Assumptions 1-2 are satisfied,*
 176 *$D^* \in \mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$ and $a \in (0, 1]$, and \mathcal{F}_{FNN} is the function class*
 177 *of ReLU DNNs with width \mathcal{W} and depth \mathcal{D} specified in (9). Then, for $M \geq 1$ and $n \geq \text{Pdim}(\mathcal{F}_{\text{FNN}})$,*
 178 *we have*

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \leq C \left(\xi_n^2 + C_1 (KL)^{-\frac{4\beta}{d}} \right),$$

179 where $C_1 = (\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)}$ and the constant C depends only on (μ, σ, M) .

180 Furthermore, if

$$\mathcal{W} = 114(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1}, \mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 \left[n^{\frac{d}{2(d+2\beta)}} \log_2 \left(8n^{\frac{d}{2(d+2\beta)}} \right) \right],$$

181 then

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \leq C_0 (\lfloor \beta \rfloor + 1)^9 d^{2\lfloor \beta \rfloor + (\beta \vee 3)} n^{-\frac{2\beta}{d+2\beta}}, \quad (10)$$

182 where the constant C_0 depends only on (μ, σ, M) .

183 The convergence rate in (10) is optimal. This can be seen by considering a density estimation problem
 184 with i.i.d observations $\{Z_{q,i}^{(1)}\}_{i=1}^{m_q}$ from an underlying unknown density q_1 on $[0, 1]^d$. To estimate
 185 q_1 , we sample referencing observations $\{Z_{p,j}^{(1)}\}_{j=1}^{m_p}$ with $m_p \geq m_q$, from a uniform distribution
 186 $\text{Unif}([0, 1]^d)$ whose density $p_1 \equiv 1$. Thus, estimating the density ratio q_1/p_1 is equivalent to
 187 estimating q_1 . According to (4), we obtain the estimator \hat{q}_1 of q_1 . If $\log q_1 \in \mathcal{H}^\beta([0, 1]^d, M)$ where
 188 $\beta = k + a$ with $k \in \mathbb{N}^+$ and $a \in (0, 1]$, a neural estimator based on the network structure specified
 189 in Theorem 3 satisfies

$$E_{p_1, q_1} \|\hat{q}_1 - q_1\|_{\max}^2 \leq C_0 (\lfloor \beta \rfloor + 1)^9 d^{2\lfloor \beta \rfloor + (\beta \vee 3)} m_q^{-\frac{2\beta}{d+2\beta}}. \quad (11)$$

190 Tsybakov (2008) showed that for a density belonging to the Hölder function class, the optimal
 191 minimax rate of the density estimation is $O_p \left(m_q^{-2\beta/(d+2\beta)} \right)$. Hence, our estimator achieves the
 192 optimal minimax rate.

193 In addition, the existing error bounds usually contain a prefactor depending on the dimension d
 194 exponentially, e.g. 2^d (Devroye & Lugosi, 1996). Such a prefactor can be very large even for a
 195 moderately large d , which severely degrades the quality of an error bound. The prefactors in our
 196 results depend on d only polynomially and are much smaller than those in the existing bounds.

197 Under Assumption 2, to derive a nonasymptotic error bound for the log-density ratio estimator \widehat{R} , we
 198 note that

$$E_{p^*, q^*} \|\widehat{R} - R^*\|_{\max}^2 \leq e^{2M} E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2.$$

199 Thus a bound for \widehat{R} follows directly from a bound for \widehat{D} .

200 **Remark 2.** *Appendix A.2 contains some examples of p^* and q^* such that $D^* = \log(q^*/p^*) \in$*
 201 *$\mathcal{H}^\beta([0, 1]^d, M)$.*

202 **Remark 3.** *The hypercube $[0, 1]^d$ assumption for the density ratio is made for technical*
 203 *convenience. With an unbounded support, we can bound $\|D_{\text{NN}} - D^*\|_{\max}$ using a truncation*
 204 *technique under some suitable additional assumptions, at a small price of an additional loga-*
 205 *rithm term in the error bound. Specifically, suppose the pdfs are supported on \mathbb{R}^d . In addition*
 206 *to Assumptions 1-2 and the Hölder class assumption in Theorem 3, we need to further assume*
 207 *that $\max\{E_{p^*} I(\|Z\|_\infty \geq \log n), E_{q^*} I(\|Z\|_\infty \geq \log n)\} \leq n^{-\frac{2\beta}{d+2\beta}}$. For $I = p$ or q , and any*
 208 *$D \in \mathcal{F}_{\text{FNN}}$, where \mathcal{F}_{FNN} is the function class of ReLU FNNs with width \mathcal{W} and depth \mathcal{D} specified by*

$$\mathcal{W} = 114(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1}, \mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 \left[n^{\frac{d}{2(d+2\beta)}} \log_2 \left(8n^{\frac{d}{2(d+2\beta)}} \right) \right],$$

209 we have

$$\|D_{\text{NN}} - D^*\|_{\max}^2 \leq 328M^2 (\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)} (2 \log n)^{2\lfloor \beta \rfloor} n^{-\frac{2\beta}{d+2\beta}}.$$

210 Compared with the upper bound of the approximation error in Theorem 3, when the pdfs are supported
 211 on \mathbb{R}^d (unbounded case), we can derive a similar approximation error upper bound with an additional
 212 logarithmic factor $(2 \log n)^{2\lfloor \beta \rfloor}$. The details are given in Appendix A.3.

213 3.3 Circumventing the curse of dimensionality

214 In many modern statistical and machine learning tasks, such as image processing and text analysis,
 215 the dimensionality d of the data can be high, which results in a very slow convergence rate even with
 216 a large sample size. This is known as the curse of dimensionality. Nonetheless, the data in various
 217 applications has been demonstrated to be supported or approximately supported in some subspaces
 218 or subsets with low intrinsic dimensionality (Nakada & Imaizumi, 2020). For regression problems,
 219 Nakada & Imaizumi (2020) have shown that DNNs can adaptively estimate the regression function
 220 through the low-dimensional structure of the data, and the resulting convergence rates no longer
 221 depend on the nominal high dimensionality d of the data, but on its low intrinsic dimension.

222 Motivated by these advancements, we assume that the data is concentrated on an approximate compact
 223 Riemannian submanifold \mathcal{M} with the Riemannian dimension $d_{\mathcal{M}} \ll d$.

224 **Assumption 4.** *The target log-density ratio $D^* \in \mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where
 225 $k \in \mathbb{N}^+$ and $a \in (0, 1]$, and the data from the densities p^*, q^* are concentrated on a set $\mathcal{M}_\rho \subseteq [0, 1]^d$
 226 defined as*

$$\mathcal{M}_\rho := \{x \in [0, 1]^d : \text{there exists } y \in \mathcal{M}, \|x - y\|_2 \leq \rho\},$$

227 where \mathcal{M} is a compact $d_{\mathcal{M}}$ -dimensional Riemannian submanifold and $\rho \in (0, 1)$.

228 **Theorem 4.** *Suppose Assumptions 1, 2 and 4 hold. Suppose that $D^* \in \mathcal{H}^\beta([0, 1]^d, M)$ with
 229 $\beta = k + a$, $k \in \mathbb{N}^+$ and $a \in (0, 1]$. If \mathcal{F}_{FNN} is the function class of ReLU FNNs with width and depth*

$$\mathcal{W} = 38(\lfloor \beta \rfloor + 1)^2 d_\delta^{\lfloor \beta \rfloor + 1} L \lceil \log_2(8L) \rceil, \mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 K \lceil \log_2(8K) \rceil,$$

230 where $K, L \in \mathbb{N}^+$ and $d_\delta = O(d_{\mathcal{M}} \log(d/\delta)/\delta^2) \ll d$, then when $M \geq 1$, $n > \text{Pdim}(\mathcal{F}_{\text{FNN}})$ and

$$\rho \leq (\lfloor \beta \rfloor + 1)^2 2^\beta d^{\beta - \frac{1}{2}} d_\delta^{\lfloor \beta \rfloor + (\beta - 1/2) \vee (1/2)} (KL)^{-\frac{2\beta}{d_\delta}},$$

231 we have

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \leq C(1 - \delta)^{-2\beta} \left[\xi_n^2 + C_2 (KL)^{-\frac{4\beta}{d_\delta}} \right],$$

232 where the constant C only depends on (μ, σ, M) , $C_2 = (\lfloor \beta \rfloor + 1)^4 (2d)^{2\beta} d_\delta^{3\beta + (\beta \vee 1)}$, and ξ_n is
 233 defined in (8).

234 By Theorem 4, if we set $\mathcal{W} = 114(\lfloor \beta \rfloor + 1)^2 d_\delta^{\lfloor \beta \rfloor + 1}$, $\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 \lceil n^{\zeta_\delta} \log_2(8n^{\zeta_\delta}) \rceil$, with
 235 $\zeta_\delta = d_\delta / (2(d_\delta + 2\beta))$, then

$$E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \leq C_0 C_3 (1 - \delta)^{-2\beta} n^{-\frac{2\beta}{d_\delta + 2\beta}}, \quad (12)$$

236 where the constant C_0 only depends on (μ, σ, M) and $C_3 = (\lfloor \beta \rfloor +$
 237 $1)^9 \max\{d_\delta^{2\lfloor \beta \rfloor + 3}, (2d)^{2\beta} d_\delta^{3\beta + (\beta \vee 1)}\}$. The convergence rate $n^{-2\beta/(d_\delta + 2\beta)}$ in (12) only de-
 238 pends on $d_\delta \ll d$, instead of the ambient dimension d . Therefore, Theorem 4 shows that a
 239 low-dimensional Riemannian manifold support assumption can alleviate the curse of dimensionality.

240 4 Error analysis of the telescoping density-ratio estimator

241 When the difference or the ‘gap’ between two densities is large, a single-ratio estimation method
 242 may perform poorly. This is referred to as the *density-chasm problem* (Rhodes et al., 2020). To
 243 alleviate this problem, Rhodes et al. (2020) proposed an approach called Telescoping density-Ratio
 244 Estimation (TRE). This approach first gradually transports samples from q^* to samples from p^* ,
 245 creating a chain of intermediate datasets, then estimates the density ratio between consecutive datasets
 246 along this chain. The chained ratios are combined via a telescoping product which yields an estimate
 247 of the original density ratio. The experiments conducted by Rhodes et al. (2020) demonstrate that
 248 TRE can yield substantial improvements over existing single-ratio methods for mutual information
 249 estimation, representation learning and energy-based modelling.

250 **We now provide a theoretical analysis of TRE, which partially explains why TRE performs well.** For
 251 notational simplicity, suppose $n_p = n_q \equiv n$ below.

252 For $k = 0, 1, \dots, K$, Rhodes et al. (2020) constructed a chain of intermediate samples connecting
 253 q^* and p^* by setting $Z_{k,i} = (1 - \alpha_k^2)^{1/2} Z_{q,i} + \alpha_k Z_{p,i}$, $i = 1, \dots, n$, where $0 = \alpha_0 < \alpha_1 < \dots <$
 254 $\alpha_{K-1} < \alpha_K = 1$, and used these samples to build a TRE.

255 To simplify the analysis, we use a slightly different chain of intermediate samples as follows. For
 256 $k = 0, 1, \dots, K$, let

$$Z_{k,i} = (1 - \delta_{k,i})Z_{q,i} + \delta_{k,i}Z_{p,i}, \quad i = 1, \dots, n, \quad (13)$$

257 where $\delta_{k,i}$, $i = 1, \dots, n$, are i.i.d. Bernoulli random variables with success probability α_k .

258 Let q_k be the density of the synthetic data $Z_{k,i}$ constructed this way. We have $q_k(z) = (1 - \alpha_k)q^*(z) +$
 259 $\alpha_k p^*(z)$, $k = 1, \dots, K - 1$. Therefore, the distribution of the samples from q_k in the chain is a
 260 simple mixture of q^* and p^* with the mixing proportions $1 - \alpha_k$ and α_k , instead of a more complex
 261 convolution of two densities using the construction of Rhodes et al. (2020). As α_k changes from
 262 $\alpha_0 = 0$ to $\alpha_K = 1$ over a grid $\{\alpha_0, \alpha_1, \dots, \alpha_K\} \subset [0, 1]$, the distributions of the samples in the
 263 chain move gradually from q^* to p^* . Let $q_0 = q^*$ and $q_K \equiv p^*$. Then,

$$R^*(z) = \frac{q^*(z)}{p^*(z)} = \prod_{i=0}^{K-1} R_i^*(z), \quad z \in \mathcal{Z}, \quad (14)$$

264 where $R_i^*(z) = q_i(z)/q_{i+1}(z)$. For $k = 0, 1, \dots, K - 1$, applying the neural density-ratio estimator
 265 with $\{Z_{k,j}\}_{j=1}^{n_k}$ and $\{Z_{k+1,j}\}_{j=1}^{n_{k+1}}$ yields an estimator \widehat{R}_i^* of R_i^* . Then the telescoping density ratio
 266 estimator of R^* is $\prod_{i=0}^{K-1} \widehat{R}_i^*$.

267 We consider the log-density ratio. Let \widehat{D}_k be the neural estimator of $D_k^* \equiv \log(q_k/q_{k+1})$. Based on
 268 (14), the telescoping estimator of the log-density ratio $D^* \equiv \log R^*$ is

$$\widehat{D}_{\text{TRE}} = \sum_{k=0}^{K-1} \widehat{D}_k. \quad (15)$$

269 In what follows, we show that under certain conditions, the telescoping estimator has an improved
 270 asymptotic error bound. The intuition is as follows: when q_k/q_{k+1} is bounded or $q_k(z)/q_{k+1}(z) \ll$
 271 $q^*(z)/p^*(z)$ for $z \in \mathcal{Z}$, where q_k and q_{k+1} are the densities of the synthetic data $\{Z_{k,j}\}_{j=1}^n$ and
 272 $\{Z_{k+1,j}\}_{j=1}^n$, respectively, the truncation error for q_k/q_{k+1} vanishes or is far less than that for q^*/p^* .
 273 This can help the telescoping density-ratio estimator perform better than a single-ratio estimator.

274 Assume that $q^* \geq c_1$ and $c_1 \leq p^* \leq c_2$, where $0 < c_1, c_2 < \infty$ are two constants. Thus,
 275 $D^* = \log(q^*/p^*) \geq \log(c_1/c_2)$. Therefore, Assumption 3 is satisfied. For any finite set $\mathcal{A} \subset \mathbb{R}$,
 276 $\max \mathcal{A}$ denotes the maximal value in \mathcal{A} . Let $M = \log(c_2/c_1)$ and M_0 be a constant satisfying

$$M_0 \geq \max_{M, \alpha} \mathcal{A}_{M, \alpha}^{(2K)}, \quad (16)$$

277 where

$$\mathcal{A}_{M, \alpha}^{(2K)} = \left\{ M, 1 \right\} \cup \left\{ \log \frac{1 - \alpha_{k-1}}{1 - \alpha_k}, 1 \leq k \leq K - 1 \right\} \cup \left\{ \log \frac{(e^M - 1)\alpha_k + 1}{(e^M - 1)\alpha_{k-1} + 1}, 1 \leq k \leq K - 1 \right\}.$$

278 Based on Theorem 2, we can establish an asymptotic error bound for the telescoping estimator \widehat{D}_{TRE}
 279 defined in (15), with

$$\widehat{D}_k \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}^0} \widehat{\mathcal{B}}_{\psi}^k(e^D),$$

where

$$\widehat{\mathcal{B}}_{\psi}^k(e^D) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_1(D(Z_{k+1,i})) + \frac{1}{n} \sum_{i=1}^n \mathcal{L}_2(D(Z_{k,i})),$$

280 where \mathcal{L}_1 and \mathcal{L}_2 are defined in (5) and $\mathcal{F}_{\text{FNN}}^0$ consists of DNNs D with $\|D\|_{\infty} \leq M_0$. To
 281 demonstrate the advantages of the telescoping estimator, we also consider the single-ratio estimator
 282 (SRE), $\widehat{D}_{\text{SRE}} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}^0} \widehat{\mathcal{B}}_{\psi}(e^D)$.

Proposition 1. Assume that $q^* \geq c_1$, $c_1 \leq p^* \leq c_2$, where the constants $0 < c_1 \leq c_2 < \infty$, and the samples $\{Z_{q,i}\}_{i=1}^n$ from q^* and $\{Z_{p,j}\}_{j=1}^n$ from p^* are independent. Then, there exists a constant $C_0(\mu, \sigma, c_1)$ depending only on (μ, σ, c_1) such that for

$$B_{\text{SRE}} = e^{M_0} C_0(\mu, \sigma, c_1) \|R^* - R_{M_0}^*\|_p,$$

283 we have

$$\begin{aligned} \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{\text{SRE}} - D^*\|_2 &\leq B_{\text{SRE}}, \\ \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{\text{TRE}} - D^*\|_2 &\leq (1 - \alpha_{K-1}) B_{\text{SRE}}, \end{aligned}$$

284 where $\|f\|_2 = [\int_{\mathcal{Z}} f^2(z) dz]^{1/2}$ for any square integrable function f .

285 Proposition 1 shows that for a given sequence $0 = \alpha_0 < \alpha_1 < \dots < \alpha_{K-1} < \alpha_K = 1$ and a
286 truncation level M_0 , the upper bound for the asymptotic L_2 -error of \widehat{D}_{TRE} is reduced by a factor
287 $(1 - \alpha_{K-1})$ with $0 < 1 - \alpha_{K-1} < 1$. This upper bound can be far less than that of \widehat{D}_{SRE} when
288 α_{K-1} is close to 1. Therefore, TRE can improve the asymptotic error bound over the bound for the
289 single-ratio method.

290 It is important to note that there is a tradeoff between the value of α_{K-1} and the truncation level
291 M_0 dictated by (16). For instance, with $\alpha_1 \leq \dots \leq \alpha_{K-2}$ fixed, the closer α_{K-1} is to 1, in
292 view of (16), the larger M_0 is. Larger α_{K-1} sharpens the pre-factor $(1 - \alpha_{K-1})$ and larger M_0 also
293 improves $\|R^* - R_{M_0}^*\|_p$, but deteriorates the pre-factor e^{M_0} .

294 Proposition 1 is generally not applicable to the original chain of TRE. The difficulty is due to
295 the possibly intensive oscillation of density ratios caused by the convolution form for the density
296 of the sum of two random variables. We illustrate this by a toy example: suppose Z_q, Z_p are
297 i.i.d. uniform random variables on $[0, 1]$. For any $t \in (0, 1/2]$, $(1-t)Z_q + tZ_p$ has density
298 $q_t^*(z) = \frac{z}{t(1-t)} I\{0 \leq z \leq t\} + \frac{1}{1-t} I\{t < z \leq 1-t\} + \frac{1-z}{t(1-t)} I\{1-t < z \leq 1\}$. In this case, q^*/q_t^*
299 is unbounded and oscillates sharply when z is close to 0 or 1. This makes it hard to estimate q^*/q_t^* .
300 However, the chain we used does not have this problem, which may be a good choice in practice.

301 Additionally, we conduct simulation studies to evaluate the performance of our proposed mixing
302 chain and the original convolution chain; see Table 2 for the results. The simulation settings are
303 given in Appendix A.4. Table 2 shows that, for the models considered in the simulation studies, the
304 proposed mixing chain performs comparably or better compared with the original convolution chain.

Table 2: The MSEs averaged over 10 replications and the corresponding standard errors in parentheses between the telescoping ratio estimate (TRE) of log density-ratio and its true value for the proposed mixing chain (mTRE) and the original convolution chain (cTRE) under different settings, where n is the training data sample size and K is the length of the chain. The bold one is the best in a specific setting among the two estimates.

Setting	Method	(n,K)			
		(5000,5)	(5000,10)	(10000,5)	(10000,10)
Beta	mTRE(ours)	0.9850(0.0269)	0.8840(0.0180)	1.0109(0.0171)	0.9299(0.0194)
	cTRE	1.4670(0.0606)	1.2935(0.0274)	1.3674(0.0625)	1.2850(0.0293)
Normal	mTRE(ours)	2.7426(0.0370)	2.8330(0.0450)	2.7483(0.0367)	2.7813(0.0265)
	cTRE	2.7987(0.0586)	2.7076(0.0293)	2.8184(0.0347)	2.7503(0.0297)

305 5 Related work: comparison with the NN-BD estimator

306 There has been much work on the error analysis of nonparametric density-ratio estimation (Nguyen
307 et al., 2010; Sugiyama et al., 2008; Kanamori et al., 2012a; Yamada et al., 2013). These results show
308 that when the targeted density-ratio belongs to certain function space \mathcal{H} , such as RKHS, and thus
309 no approximation error is incurred, their estimators achieve certain nonparametric convergence rate

310 decided by the complexity of \mathcal{H} . Compared to these works, our results consider the approximation
 311 error using neural network functions and still achieves the minimax optimality under some mild
 312 conditions.

313 Our work is most related to the paper by Kato & Teshima (2021), who proposed a non-negative
 314 Bregman divergence (NN-BD) method to tackle the possible over-fitting problem due to the unbound-
 315 edness of certain Bregman divergences. We compare our theoretical results with those for the NN-BD
 316 estimator of Kato & Teshima (2021). Using the notation in this paper, we restate two conditions
 317 required in Kato & Teshima (2021):

318 (a) Let $\mathcal{F}_{\text{FNN}}^R$ be a class of FNNs with output taking values in $[e^{-M}, e^M]$ for some finite $M > 0$.
 319 The target density-ratio $R^* \in \mathcal{F}_{\text{FNN}}^R$. Moreover, for any function in $\mathcal{F}_{\text{FNN}}^R$, its Frobenius norm of
 320 the parameter matrix W_j in the j th layer is bounded by $\mathcal{B}_j \geq 0$ and the activation functions are
 321 1-Lipschitz positive-homogeneous.

322 (b) The function $\psi(\cdot)$ is σ -strongly convex. Let $\ell_1(t) = \psi^*(t)t - \psi(t) + A$, $\ell_2(t) = -\tilde{\psi}(t)$, $t \in$
 323 $[e^{-M}, e^M]$, where $\psi^*(t) = C_{nn}\{\psi'(t)t - \psi(t)\} + \tilde{\psi}(t)$. Here $\tilde{\psi}(t)$ is a function bounded above, C_{nn}
 324 and A are user-selected constants. Suppose $\ell_1(\cdot)$ and $\ell_2(\cdot)$ are Lipschitz functions on $[e^{-M}, e^M]$.

325 Under these two conditions, Kato & Teshima (2021) rewrote the BD in (3) as

$$\mathcal{B}_\psi(R) = E_p \ell_1(R(Z)) - C_{nn} E_q \ell_1(R(Z)) + E_q \ell_2(R(Z)) + (1 - C_{nn})A, \quad (17)$$

326 and proposed the density-ratio estimator \hat{R}_{KT} defined as

$$\hat{R}_{\text{KT}} \in \arg \min_{R \in \mathcal{F}_{\text{FNN}}} \left\{ \frac{1}{n_q} \sum_{i=1}^{n_q} \ell_2(R(Z_{q,i})) + \left[\frac{1}{n_p} \sum_{j=1}^{n_p} \ell_1(R(Z_{p,j})) - \frac{C_{nn}}{n_q} \sum_{j=1}^{n_q} \ell_1(R(Z_{q,j})) \right]_+ \right\},$$

327 where $[a]_+ = \max(0, a)$ for any $a \in \mathbb{R}$. They showed that

$$\|\hat{R}_{\text{KT}} - R^*\|_p = O_p\left(n^{-1/(2+a)}\right) \quad (18)$$

328 for any $0 \leq a \leq 2$.

329 According to Theorem 1, we have the following corollary for our density-ratio estimator $\hat{R} = \exp(\hat{D})$.

330 **Corollary 2.** *Under Assumption 1, when $n \geq \text{Pdim}(\mathcal{F}_{\text{FNN}}^R)$, there exists a constant C depending*
 331 *only on (μ, σ, M) such that, for any $\gamma \geq 0$, with probability at least $1 - \exp(-\gamma)$,*

$$\|\hat{R} - R^*\|_p \leq C \left(\sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}^R) \log n}{n}} + \sqrt{\frac{\gamma}{n}} \right).$$

332 Corollary 2 implies that $\|\hat{R} - R^*\|_p = O_p(\sqrt{\log n/n})$, when the true density-ratio $R^* \in \mathcal{F}_{\text{FNN}}$. This
 333 convergence rate is slightly faster than the rate for \hat{R}_{KT} given in (18). Moreover, the boundedness
 334 assumption for the weights of the neural network functions, as imposed by Kato & Teshima (2021),
 335 is not needed. Corollary 2 also shows that, if the target ratio is assumed to belong to the optimization
 336 space (or hypothesis space), i.e., $R^* \in \mathcal{F}_{\text{FNN}}^R$ without approximation error, then the convergence rate
 337 does not depend on the dimension of the data. In other words, the estimation of R^* does not suffer
 338 from the curse of dimensionality. However, this is probably not realistic. Therefore, it is important to
 339 consider the approximation error due to the fact that $R^* \notin \mathcal{F}_{\text{FNN}}^R$ in applications.

340 6 Conclusions

341 In this paper, we have established the non-asymptotic error bounds for the deep density-ratio estimator
 342 using the Bregman divergence criterion. Under reasonable conditions, we have shown that the deep
 343 density-ratio estimator achieves the optimal minimax convergence rate. When the data is supported
 344 on an approximate low-dimensional manifold, we have shown that the neural estimator can mitigate
 345 the curse of dimensionality. We have also analyzed the convergence properties of the telescoping
 346 density ratio estimator (Rhodes et al., 2020) and provided sufficient conditions under which it has a
 347 lower error bound than a single-ratio estimator.

348 A limitation of this work is that certain boundedness assumptions on the target density ratio such
349 as Assumption 2 or 3 is needed. Also, when the boundedness assumption is partially relaxed as in
350 Assumption 3, the error bound in Theorem 2 is not as sharp as that with the boundedness assumption
351 in Theorem 1. It would be interesting to further relax or remove such assumptions. It would also be
352 useful to improve the error bound in Theorem 2 if possible. These are interesting and challenging
353 problems that deserve further study in the future.

354 **References**

- 355 Anthony, M. and Bartlett, P. L. *Neural Network Learning: Theoretical Foundations*. Cambridge
356 University Press, Cambridge, 1999.
- 357 Baraniuk, R. G. and Wakin, M. B. Random projections of smooth manifolds. *Found. Comput. Math.*,
358 9(1):51–77, 2009.
- 359 Bartlett, P. L., Bousquet, O., and Mendelson, S. Local rademacher complexities. *The Annals of*
360 *Statistics*, 33(4):1497–1537, 2005.
- 361 Bartlett, P. L., Harvey, N., Liaw, C., and Mehrabian, A. Nearly-tight VC-dimension and pseudodi-
362 mension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20
363 (63):1 – 17, 2019.
- 364 Beck, A. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics,
365 2017.
- 366 Bregman, L. The relaxation method of finding the common point of convex sets and its application
367 to the solution of problems in convex programming. *USSR Computational Mathematics and*
368 *Mathematical Physics*, 7(3):200 – 217, 1967.
- 369 Ceylan, C. and Gutmann, M. U. Conditional Noise-Contrastive Estimation of Unnormalised Models.
370 In *Proceedings of the 35th International Conference on Machine Learning*, pp. 725–733. PMLR,
371 2018.
- 372 Cox, T. F. and Ferry, G. Robust logistic discrimination. *Biometrika*, 78(4):841 – 849, 1991.
- 373 Devroye, L. and Lugosi, G. A universally acceptable smoothing factor for kernel density estimates.
374 *The Annals of Statistics*, 24(6):2499 – 2512, 1996.
- 375 Fefferman, C. Whitney’s extension problem for c^m . *Annals of Mathematics.*, 164(1):313–359, 2006.
- 376 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and
377 Bengio, Y. Generative adversarial nets. In *Advances in neural information processing systems*, pp.
378 2672–2680, 2014.
- 379 Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., and Schölkopf, B. Covariate shift
380 by kernel mean matching. In Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N.
381 (eds.), *Dataset Shift in Machine Learning*, chapter 8, pp. 131 – 160. MIT Press, Cambridge, 2009.
- 382 Gutmann, M. and Hyvärinen, A. Noise-contrastive estimation of unnormalized statistical models,
383 with applications to natural image statistics. *Journal of Machine Learning Research*, 13:307–361,
384 2012.
- 385 Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional
386 networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.
387 2261–2269, 2017.
- 388 Jiao, Y., Shen, G., Lin, Y., and Huang, J. Deep nonparametric regression on approximately low-
389 dimensional manifolds. *arXiv:2104.06708*, 2021.
- 390 Kanamori, T., Hido, S., and Sugiyama, M. A least-squares approach to direct importance estimation.
391 *Journal of Machine Learning Research*, 10(48):1391 – 1445, 2009.
- 392 Kanamori, T., Suzuki, T., and Sugiyama, M. Statistical analysis of kernel-based least-squares
393 density-ratio estimation. *Machine Learning*, 86(3):335 – 367, 2012a.
- 394 Kanamori, T., Suzuki, T., and Sugiyama, M. f -divergence estimation and two-sample homogeneity
395 test under semiparametric density-ratio models. *IEEE Transactions on Information Theory*, 58:708
396 – 720, 2012b.
- 397 Kato, M. and Teshima, T. Non-negative bregman divergence minimization for deep direct density
398 ratio estimation. In *International Conference on Machine Learning*, pp. 5320 – 5333, 2021.

- 399 Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- 400 LeCun, Y., Cortes, C., and Burges, C. Mnist handwritten digit database. *AT&T Labs [Online]*.
401 Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- 402 Moustakides, G. V. and Basioti, K. Training neural networks for likelihood/density ratio estimation.
403 *arXiv:1911.00405*, 2019.
- 404 Nakada, R. and Imaizumi, M. Adaptive approximation and estimation of deep neural network with
405 intrinsic dimensionality. *Journal of Machine Learning Research*, 21(174):1 – 38, 2020.
- 406 Nguyen, X., Wainwright, M., and Jordan, M. Estimating divergence functionals and the likelihood
407 ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847 – 5861,
408 2010.
- 409 Nowozin, S., Cseke, B., and Tomioka, R. f-gan: Training generative neural samplers using variational
410 divergence minimization. In *Advances in neural information processing systems*, pp. 271–279,
411 2016.
- 412 Parzen, E. On Estimation of a Probability Density Function and Mode. *The Annals of Mathematical*
413 *Statistics*, 33(3):1065 – 1076, 1962.
- 414 Qin, J. Inferences for case-control and semiparametric two-sample density ratio models. *Biometrika*,
415 85(3):619 – 630, 1998.
- 416 Rhodes, B., Xu, K., and Gutmann, M. U. Telescoping density-ratio estimation. In *Advances in Neural*
417 *Information Processing Systems*, 2020.
- 418 Rosenblatt, M. Remarks on some nonparametric estimates of a density function. *The Annals of*
419 *Mathematical Statistics*, 27(3):832 – 837, 1956.
- 420 Silverman, B. W. Density ratios, empirical likelihood and cot death. *Journal of the Royal Statistical*
421 *Society. Series C (Applied Statistics)*, 27(1):26 – 33, 1978.
- 422 Sugiyama, M., Nakajima, S., Kashima, H., Bunau, P. V., and Kawanabe, M. Direct importance
423 estimation for covariate shift adaptation. *Annals of the Institute of Statistical Mathematics*, 60(4):
424 699 – 746, 2008.
- 425 Sugiyama, M., Suzuki, T., Itoh, Y., Kanamori, T., and Kimura, M. Least-squares two-sample test.
426 *Neural Networks*, 24(7):735 – 751, 2011.
- 427 Sugiyama, M., Suzuki, T., and Kanamori, T. Density-ratio matching under the bregman divergence:
428 a unified framework of density-ratio estimation. *Annals of the Institute of Statistical Mathematics*,
429 64(5):1009 – 1044, 2012a.
- 430 Sugiyama, M., Suzuki, T., and Kanamori, T. *Density Ratio Estimation in Machine Learning*.
431 Cambridge University Press, Cambridge, 2012b.
- 432 Tsuboi, Y., Kashima, H., Hido, S., Bickel, S., and Sugiyama, M. Direct density ratio estimation for
433 large-scale covariate shift adaptation. *Journal of Information Processing*, 17(2):138 – 155, 2009.
- 434 Tsybakov, A. B. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer
435 New York, 2008.
- 436 Yamada, M. and Sugiyama, M. Direct importance estimation with gaussian mixture models. *IEICE*
437 *Transactions on Information and Systems*, E92.D(10):2159 – 2162, 2009.
- 438 Yamada, M., Sugiyama, M., Wichern, G., and Simm, J. Direct importance estimation with a mixture
439 of probabilistic principal component analyzers. *IEICE Transactions on Information and Systems*,
440 E93.D(10):2846 – 2849, 2010.
- 441 Yamada, M., Suzuki, T., Kanamori, T., Hachiya, H., and Sugiyama, M. Relative density-ratio
442 estimation for robust distribution comparison. *Neural computation*, 25(5):1324–1370, 2013.

443 **Checklist**

- 444 1. For all authors...
- 445 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
446 contributions and scope? [Yes]
- 447 (b) Did you describe the limitations of your work? [Yes] See Section 6.
- 448 (c) Did you discuss any potential negative societal impacts of your work? [No] There is
449 no such a potential negative societal impact in our work.
- 450 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
451 them? [Yes]
- 452 2. If you are including theoretical results...
- 453 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assump-
454 tions 1, 2, 3 and 4.
- 455 (b) Did you include complete proofs of all theoretical results? [Yes] See the subsection
456 A.1 in Appendix A.
- 457 3. If you ran experiments...
- 458 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
459 mental results (either in the supplemental material or as a URL)? [Yes] .
- 460 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
461 were chosen)? [Yes] .
- 462 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
463 ments multiple times)? [Yes] .
- 464 (d) Did you include the total amount of compute and the type of resources used (e.g., type
465 of GPUs, internal cluster, or cloud provider)? [Yes] .
- 466 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 467 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 468 (b) Did you mention the license of the assets? [N/A]
- 469 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 470
- 471 (d) Did you discuss whether and how consent was obtained from people whose data you’re
472 using/curating? [N/A]
- 473 (e) Did you discuss whether the data you are using/curating contains personally identifiable
474 information or offensive content? [N/A]
- 475 5. If you used crowdsourcing or conducted research with human subjects...
- 476 (a) Did you include the full text of instructions given to participants and screenshots, if
477 applicable? [N/A]
- 478 (b) Did you describe any potential participant risks, with links to Institutional Review
479 Board (IRB) approvals, if applicable? [N/A]
- 480 (c) Did you include the estimated hourly wage paid to participants and the total amount
481 spent on participant compensation? [N/A]

482 **A Appendix**

483 **A.1 Theoretical Proofs**

484 In the appendix, we provide all the technical details and proofs of the theorems stated in the paper.

485 **Verification of (1):** Equation (1 holds because

$$\begin{aligned} & E_{p^*} \Delta_{\psi}(R^*(Z), R(Z)) \\ &= E_{p^*} [\psi(R^*(Z)) - \psi(R(Z)) - \psi'(R(Z))(R^*(Z) - R(Z))] \\ &= E_{p^*} [\psi'(R(Z))R(Z) - \psi(R(Z))] - E_{p^*} [\psi'(R(Z))R^*(Z)] + E_{p^*} [\psi(R^*(Z))] \\ &= E_{p^*} [\psi'(R(Z))R(Z) - \psi(R(Z))] - E_{q^*} [\psi'(R(Z))] + E_{p^*} [\psi(R^*(Z))], \end{aligned}$$

486 where $E_{p^*}[\psi'(R(Z))R^*(Z)] = E_{q^*}[\psi'(R(Z))]$ by the definition of R^* . This verifies (1).

487 We now prove the following lemmas.

488 **Lemma A.1.** 1. If the convex function $f : \mathbb{R} \rightarrow \mathbb{R}$ is μ -smooth over \mathbb{R} , then for any
489 $x, y \in \mathbb{R}$, the following inequality holds

$$f(y) \leq f(x) + f'(x)(y - x) + \frac{\mu}{2}(y - x)^2.$$

490 2. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a first-order differentiable and convex function. If f is σ -strongly convex,
491 then for any $x, y \in \mathbb{R}$, the following inequality holds

$$f(y) \geq f(x) + f'(x)(y - x) + \frac{\sigma}{2}(y - x)^2.$$

492 *Proof of Lemma A.1.* The proof of Lemma A.1 is standard and can be found in Beck (2017). \square

493 **Lemma A.2.** Under Assumptions 1-2, we have

494 (a). There exist two constants $c_0 = \frac{\sigma e^{-3M}}{2}$, $C_0 = \frac{\mu e^{3M}}{2}$, such that

$$c_0 \|D - D^*\|_{\max}^2 \leq \mathcal{B}_\psi(e^D) - \mathcal{B}_\psi(e^{D^*}),$$

495 and

$$\mathcal{B}_\psi(e^D) - \mathcal{B}_\psi(e^{D^*}) \leq C_0 \|D - D^*\|_{\max}^2.$$

496 (b). For $t_1, t_2 \in [-M, M]$, there exist two constants C_1, C_2 , such that

$$|\mathcal{L}_1(t_1) - \mathcal{L}_1(t_2)| \leq C_1 |t_1 - t_2|,$$

497 and

$$|\mathcal{L}_2(t_1) - \mathcal{L}_2(t_2)| \leq C_2 |t_1 - t_2|.$$

498 Actually, we can take $C_1 = 2e^{2M}\mu$, $C_2 = e^M\mu$.

499 *Proof of Lemma A.2.* (a) Let $\Delta_\psi(x, y) := \psi(x) - \psi(y) - \psi'(x)(x - y)$. Since
500 $E_{p^*}\Delta_\psi(e^{D(Z)}, e^{D^*(Z)}) = \mathcal{B}_\psi(e^D) - \mathcal{B}_\psi(e^{D^*})$ and ψ is μ -smooth and σ -strongly convex, by Lemma
501 A.1,

$$\frac{\sigma}{2} E_{p^*} \{e^{D(Z)} - e^{D^*(Z)}\}^2 \leq E_{p^*} \Delta_\psi(e^{D(Z)}, e^{D^*(Z)}) \leq \frac{\mu}{2} E_{p^*} \{e^{D(Z)} - e^{D^*(Z)}\}^2,$$

502 and then by Assumption 2,

$$\frac{\sigma e^{-2M}}{2} E_{p^*} \{D(Z) - D^*(Z)\}^2 \leq E_{p^*} \Delta_\psi(e^{D(Z)}, e^{D^*(Z)}) \leq \frac{\mu e^{2M}}{2} E_{p^*} \{D(Z) - D^*(Z)\}^2. \quad (\text{A.1})$$

503 As $E_{p^*} \{D(Z) - D^*(Z)\}^2 = E_{q^*} e^{-D^*(Z)} \{D(Z) - D^*(Z)\}^2$ and $\|D^*\|_\infty \leq M$, we have

$$e^{-M} E_{q^*} \{D(Z) - D^*(Z)\}^2 \leq E_{p^*} \{D(Z) - D^*(Z)\}^2 \leq e^M E_{q^*} \{D(Z) - D^*(Z)\}^2. \quad (\text{A.2})$$

504 Let $c_0 = \frac{\sigma e^{-3M}}{2}$, $C_0 = \frac{\mu e^{3M}}{2}$, then (A.1) and (A.2) imply that

$$c_0 \|D - D^*\|_{\max}^2 \leq \mathcal{B}_\psi(e^D) - \mathcal{B}_\psi(e^{D^*}) \leq C_0 \|D - D^*\|_{\max}^2.$$

505 (b) Obviously, for $t_1, t_2 \in [-M, M]$,

$$\begin{aligned} |\mathcal{L}_1(t_1) - \mathcal{L}_1(t_2)| &= |\psi'(e^{t_1})e^{t_1} - \psi(e^{t_1}) - (\psi'(e^{t_2})e^{t_2} - \psi(e^{t_2}))| \\ &\leq e^{t_1} |\psi'(e^{t_1}) - \psi'(e^{t_2})| + |\psi(e^{t_1}) - \psi(e^{t_2}) - \psi'(e^{t_2})(e^{t_1} - e^{t_2})| \\ &\leq e^M \mu |e^{t_1} - e^{t_2}| + \frac{\mu}{2} |e^{t_1} - e^{t_2}|^2 \\ &\leq 2e^M \mu |e^{t_1} - e^{t_2}| \quad (\text{As } |e^{t_1} - e^{t_2}| \leq 2e^M) \\ &\leq 2e^{2M} \mu |t_1 - t_2|. \end{aligned}$$

506 and

$$\begin{aligned}
|\mathcal{L}_2(t_1) - \mathcal{L}_2(t_2)| &= |\psi'(e^{t_1}) - \psi'(e^{t_2})| \\
&\leq \mu|e^{t_1} - e^{t_2}| \\
&\leq e^M \mu|t_1 - t_2|.
\end{aligned}$$

507 The proof of the lemma is completed. \square

508 *Proof of Theorem 1.* For notational convenience, denote $\epsilon_n = \|D_{\text{NN}} - D^*\|_{\max}$ and use E_I to
509 denote E_{I^*} , $I = p, q$. Recall that E_{n_I} denotes the expectation with respect to (w.r.t) the empirical
510 distribution of $\{Z_{I,t}\}_{t=1}^{n_I}$ for $I = p, q$. As $\widehat{D} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{L}_{n_p, n_q}(D)$, where $\mathcal{L}_{n_p, n_q}(D) =$
511 $1/n_p \sum_{j=1}^{n_p} \mathcal{L}_1(D(Z_{p,j})) + 1/n_q \sum_{i=1}^{n_q} \mathcal{L}_2(D(Z_{q,i}))$, we have

$$\begin{aligned}
&c_0 \|\widehat{D} - D^*\|_{\max}^2 \\
&\leq \mathcal{B}_\psi(e^{\widehat{D}}) - \mathcal{B}_\psi(e^{D^*}) \\
&\leq \mathcal{B}_\psi(e^{\widehat{D}}) - \mathcal{B}_\psi(e^{D^*}) - \mathcal{L}_{n_p, n_q}(\widehat{D}) + \mathcal{L}_{n_p, n_q}(D_{\text{NN}}) \\
&= \mathcal{B}_\psi(e^{\widehat{D}}) - \mathcal{L}_{m, n}(\widehat{D}) - \left\{ \mathcal{B}_\psi(e^{D^*}) - \mathcal{L}_{m, n}(D^*) \right\} \\
&+ \left\{ \mathcal{L}_{n_p, n_q}(D_{\text{NN}}) - \mathcal{L}_{n_p, n_q}(D^*) \right\} \\
&= (E_{p^*} - E_{n_p})\{\mathcal{L}_1(\widehat{D}) - \mathcal{L}_1(D^*)\} + (E_q - E_{n_q})\{\mathcal{L}_2(\widehat{D}) - \mathcal{L}_2(D^*)\} \\
&+ E_{n_p}\{\mathcal{L}_1(D_{\text{NN}}) - \mathcal{L}_1(D^*)\} + E_{n_q}\{\mathcal{L}_2(D_{\text{NN}}) - \mathcal{L}_2(D^*)\}. \tag{A.3}
\end{aligned}$$

512 By Theorem 2.1 in Bartlett et al. (2005), with probability at least $1 - \exp(-\gamma_1)$,

$$E_{n_p}\{\mathcal{L}_1(D_{\text{NN}}) - \mathcal{L}_1(D^*)\} \leq E_p\{\mathcal{L}_1(D_{\text{NN}}) - \mathcal{L}_1(D^*)\} + \sqrt{2}C_1 \|D_{\text{NN}} - D^*\|_{\max} \sqrt{\frac{\gamma_1}{n}} + \frac{16C_1 M \gamma_1}{3n}. \tag{A.4}$$

513 Also, with probability at least $1 - \exp(-\gamma_1)$,

$$E_{n_q}\{\mathcal{L}_2(D_{\text{NN}}) - \mathcal{L}_2(D^*)\} \leq E_q\{\mathcal{L}_2(D_{\text{NN}}) - \mathcal{L}_2(D^*)\} + \sqrt{2}C_2 \|D_{\text{NN}} - D^*\|_{\max} \sqrt{\frac{\gamma_1}{n}} + \frac{16C_2 M \gamma_1}{3n}. \tag{A.5}$$

514 The inequalities (A.4) and (A.5) together imply that with probability at least $1 - 2 \exp(-\gamma_1)$,

$$\begin{aligned}
&E_{n_p}\{\mathcal{L}_1(D_{\text{NN}}) - \mathcal{L}_1(D^*)\} + E_{n_q}\{\mathcal{L}_2(D_{\text{NN}}) - \mathcal{L}_2(D^*)\} \\
&\leq E_p\{\mathcal{L}_1(D_{\text{NN}}) - \mathcal{L}_1(D^*)\} + E_q\{\mathcal{L}_2(D_{\text{NN}}) - \mathcal{L}_2(D^*)\} \\
&+ \sqrt{2}(C_1 + C_2) \|D_{\text{NN}} - D^*\|_{\max} \sqrt{\frac{\gamma_1}{n}} + \frac{16(C_1 + C_2)M\gamma_1}{3n} \\
&= \mathcal{B}_\psi(e^{D_{\text{NN}}}) - \mathcal{B}_\psi(e^{D^*}) + \sqrt{\frac{2\gamma_1}{n}}(C_1 + C_2) \|D_{\text{NN}} - D^*\|_{\max} + \frac{16(C_1 + C_2)M\gamma_1}{3n} \\
&\leq C_0 \|D_{\text{NN}} - D^*\|_{\max}^2 + \sqrt{\frac{2\gamma_1}{n}}(C_1 + C_2) \|D_{\text{NN}} - D^*\|_{\max} + \frac{16(C_1 + C_2)M\gamma_1}{3n}. \tag{A.6}
\end{aligned}$$

515 **Step 1.** Let $g = (D - D^*)^2$, then $g \leq 4M^2$ by Assumption 2. If $\|D - D^*\|_{\max} \leq r$, then

$$\begin{aligned}
\text{var}_p(g) \leq E_p(g^2) &= E_p(D - D^*)^4 \\
&\leq 4M^2 E_p(D - D^*)^2 \\
&\leq 4M^2 r^2.
\end{aligned}$$

516 Regarding g as a function of $D - D^*$, we have

$$\begin{aligned}
|g(D_1 - D^*) - g(D_2 - D^*)| &= |D_1^2 - 2D_1 D^* - (D_2^2 - 2D_2 D^*)| \\
&= |(D_1 + D_2 - 2D^*)(D_1 - D_2)| \\
&= |(D_1 + D_2 - 2D^*)\{(D_1 - D^*) - (D_2 - D^*)\}| \\
&\leq 4M|(D_1 - D^*) - (D_2 - D^*)|.
\end{aligned}$$

517 Thus g can be viewed as the function of $D - D^*$ with a Lipschitz constant $4M$. Denote $\mathcal{F}_{\text{FNN}}^{D^*,r} =$
518 $\{D \in \mathcal{F}_{\text{FNN}}, \|D - D^*\|_{\max} \leq r\}$, and

$$R_{n_I} \mathcal{F} = \sup_{f \in \mathcal{F}} \frac{1}{n_I} \sum_{i=1}^{n_I} \eta_i^I f(Z_{I,i}), \quad I = p, q,$$

519 where $\eta_i^I, i = 1, 2, \dots, n_I$ are i.i.d. Rademacher variables. For the rest of the proof of Theorem
520 1, we use $E_\eta R_{n_I} \mathcal{F}$ to denote the conditional expectation of $R_{n_I} \mathcal{F}$ w.r.t $\eta_i^I, i = 1, 2, \dots, n_I$, given
521 $Z_{I,i}, i = 1, 2, \dots, n_I$ and $E_{I,\eta} R_{n_I} \mathcal{F}$ to denote the expectation of $R_{n_I} \mathcal{F}$ jointly w.r.t $\eta_i^I, Z_{I,i}, i =$
522 $1, 2, \dots, n_I$. Again, by Theorem 2.1 in Bartlett et al. (2005), with probability at least $1 - \exp(-\gamma_1)$,

$$\begin{aligned} & \|D - D^*\|_{p,n_p}^2 - \|D - D^*\|_p^2 \\ & \leq 3E_{p,\eta} R_{n_p} \left\{ (D - D^*)^2 : D \in \mathcal{F}_{\text{FNN}}^{D^*,r} \right\} + 2\sqrt{\frac{2\gamma_1}{n}} M + \frac{16M^2}{3} \frac{\gamma_1}{n} \\ & \leq 24ME_{p,\eta} R_{n_p} \left\{ (D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,r} \right\} + 2\sqrt{\frac{2\gamma_1}{n}} Mr + \frac{16M^2}{3} \frac{\gamma_1}{n}, \end{aligned} \quad (\text{A.7})$$

523 where the last inequality follows from Talagrand's contraction theorem. Similarly, with probability at
524 least $1 - \exp(-\gamma_1)$,

$$\|D - D^*\|_{q,n_q}^2 - \|D - D^*\|_q^2 \leq 24ME_{q,\eta} R_{n_q} \left\{ (D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,r} \right\} + 2\sqrt{\frac{2\gamma_1}{n}} Mr + \frac{16M^2}{3} \frac{\gamma_1}{n}. \quad (\text{A.8})$$

525 Let

$$\frac{R_n(r)}{24M} = \max_{I \in \{p,q\}} \left\{ E_{I,\eta} R_{n_I} \left\{ (D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,r} \right\} \right\}.$$

526 When

$$r^2 \geq R_n(r), \quad r^2 \geq \frac{16M^2\gamma}{3n}, \quad (\text{A.9})$$

527 (A.7) and (A.8) indicate that with probability at least $1 - 2\exp(-\gamma_1)$,

$$\begin{aligned} \|D - D^*\|_{n_p, n_q}^2 &= \max\{\|D - D^*\|_{p, n_p}^2, \|D - D^*\|_{q, n_q}^2\} \\ &\leq \max\{\|D - D^*\|_p^2, \|D - D^*\|_q^2\} + R_n(r) + 2\sqrt{\frac{2\gamma_1}{n}} Mr + \frac{16M^2}{3} \frac{\gamma_1}{n} \\ &= \|D - D^*\|_{\max}^2 + R_n(r) + 2\sqrt{\frac{2\gamma_1}{n}} Mr + \frac{16M^2}{3} \frac{\gamma_1}{n} \\ &\leq (2r)^2. \end{aligned}$$

528 Thus, when (A.9) holds, with probability at least $1 - 2\exp(-\gamma_1)$,

$$\|D - D^*\|_{\max} \leq r \Rightarrow \|D - D^*\|_{n_p, n_q} \leq 2r. \quad (\text{A.10})$$

529 **Step 2.** Suppose $\|\hat{D} - D^*\|_{\max} \leq r_0$ and let

$$\mathcal{G}_i = \left\{ \mathcal{L}_i(D) - \mathcal{L}_i(D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*, r_0} \right\}, \quad i = 1, 2.$$

530 For each $(I, i) \in \{(p, 1), (q, 2)\}$, with probability at least $1 - 2\exp(-\gamma_1)$,

$$(E_I - E_{n_I}) \{ \mathcal{L}_i(\hat{D}) - \mathcal{L}_i(D^*) \} \leq 6E_\eta R_{n_I} \mathcal{G}_i + \sqrt{2} C_i r_0 \sqrt{\frac{\gamma_1}{n}} + \frac{46C_i M \gamma_1}{3n}. \quad (\text{A.11})$$

531 Denote $\hat{\mathcal{F}}_{\text{FNN}}^{D^*, r} = \{D \in \mathcal{F}_{\text{FNN}}, \|D - D^*\|_{n_p, n_q} \leq r\}$. By (A.10) in Step 1, when $r_0^2 \geq R_n(r_0)$ and
532 $r_0^2 \geq 16M^2\gamma_1/(3n)$, with probability at least $1 - 2\exp(-\gamma_1)$, for each $(I, i) \in \{(p, 1), (q, 2)\}$,

$$\begin{aligned} E_\eta R_{n_I} \mathcal{G}_i &\leq 2C_i E_\eta R_{n_I} \left\{ (D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*, r_0} \right\} \\ &\leq 2C_i E_\eta R_{n_I} \left\{ (D - D^*) : D \in \hat{\mathcal{F}}_{\text{FNN}}^{D^*, 2r_0} \right\}. \end{aligned}$$

533 Denote $\hat{\mathcal{F}}_I^{D^*,r} = \{D \in \mathcal{F}_{\text{FNN}}, \|D - D^*\|_{I,n_I} \leq r\}$. When $n \geq \text{Pdim}(\mathcal{F}_{\text{FNN}})$, $r_0 \geq 1/n$ and
 534 $n \geq (2eM)^2$, we have

$$E_\eta R_{n_I} \{(D - D^*) : D \in \hat{\mathcal{F}}_I^{D^*,2r_0}\} \leq 64r_0 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}}, \quad (\text{A.12})$$

535 and thus

$$E_\eta R_{n_I} \mathcal{G}_i \leq 128C_i r_0 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}}. \quad (\text{A.13})$$

536 Combining (A.3) (A.6) (A.11) and (A.13), with probability at least $1 - 8 \exp(-\gamma_1)$, we have

$$\begin{aligned} c_0 \|\hat{D} - D^*\|_{\max}^2 &\leq 768(C_1 + C_2)r_0 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} \\ &\quad + \sqrt{\frac{2\gamma_1}{n}}(C_1 + C_2)r_0 + \frac{46(C_1 + C_2)M\gamma_1}{3n} + C_0\epsilon_n^2 \\ &\quad + \sqrt{\frac{2\gamma_1}{n}}(C_1 + C_2)\epsilon_n + \frac{16(C_1 + C_2)M\gamma_1}{3n} \\ &= (C_1 + C_2)r_0 \left(768 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{2\gamma_1}{n}} \right) \\ &\quad + C_0\epsilon_n^2 + \sqrt{\frac{2\gamma_1}{n}}(C_1 + C_2)\epsilon_n + \frac{62(C_1 + C_2)M\gamma_1}{3n}. \end{aligned}$$

537 Therefore, when $\max \left\{ \sqrt{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n/n}, \epsilon_n \right\} \ll r_0$, there exists $r_1 \ll r_0$ such that $\|\hat{D} -$
 538 $D^*\|_{\max} \ll r_1$.

539 **Step 3.** Let $r_* = \inf\{r \geq 0 : R_n(s) \leq s^2, \text{ for } s \geq r\}$ and $E =$
 540 $\left\{ \|D - D^*\|_{n_p, n_q} \leq 4r_* \text{ for all } D \in \mathcal{F}_{\text{FNN}}^{D^*,2r_*} \right\}$. We intend to prove

$$r_* \leq \kappa M \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}}, \quad \kappa = 24 \times 130. \quad (\text{A.14})$$

541 When $r_* \leq 2\sqrt{3}M\sqrt{\log n/n}/3$, the inequality is trivial. When $r_* \geq 2\sqrt{3}M\sqrt{\log n/n}/3$, by the
 542 result in Step 1, $P(E) \geq 1 - 2/n$. As a result,

$$\begin{aligned} r_*^2 &\leq R_n(r_*) \\ &\leq R_n(2r_*) \\ &= 24M \max_{I \in \{p,q\}} \{E_{I,\eta} R_{n_I} \{(D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,2r_*}\}\}. \end{aligned}$$

543 For each $I \in \{p, q\}$,

$$\begin{aligned} E_{I,\eta} R_{n_I} \{(D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,2r_*}\} &= E_I E_\eta R_{n_I} \{(D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,2r_*}\} \\ &= E_I E_\eta R_{n_I} \{(D - D^*) : D \in \mathcal{F}_{\text{FNN}}^{D^*,2r_*}\} (I_E + I_{E^c}) \\ &\leq E_I E_\eta R_{n_I} \{(D - D^*) : D \in \hat{\mathcal{F}}_{\text{FNN}}^{D^*,4r_*}\} + \frac{4M}{n}. \end{aligned}$$

544 It follows from (A.12) that

$$\begin{aligned} r_*^2 &\leq 24M \left(128r_* \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \frac{4M}{n} \right) \\ &= 24M \left(128r_* \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + r_* \cdot \frac{4M}{n} \cdot \frac{1}{r_*} \right) \\ &\leq 24Mr_* \left(128 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{3}{n \log n}} \right) \\ &\leq \kappa \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} Mr_*, \end{aligned}$$

545 where $\kappa = 24 \times 130$. Thus, $r_* \leq \kappa M \sqrt{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n/n}$ and (A.14) is proved.

546 **Step 4.** Let $B_{\max}(D^*, r) = \{D \in \mathcal{F}_{\text{FNN}}, \|D - D^*\|_{\max} \leq r\}$, $\bar{r} \geq \max(\sqrt{\log n/n}, r_*)$ and
 547 $l = \lfloor \log_2(2M/\sqrt{\log n/n}) \rfloor$. Then, the neural network function space \mathcal{F}_{FNN} can be divided into

$$B_{\max}(D^*, \bar{r}), B_{\max}(D^*, 2\bar{r}) \setminus B_{\max}(D^*, \bar{r}), \dots, B_{\max}(D^*, 2^l \bar{r}) \setminus B_{\max}(D^*, 2^{l-1} \bar{r}).$$

548 As $\bar{r} \geq r_*$, it then follows from the definition of r_* that $\bar{r}^2 \geq R_n(\bar{r})$. Further, if $\bar{r}^2 \geq$
 549 $16M^2\gamma_1/(3n)$, according to (A.10) in Step 1, with probability at least $1 - 2l \exp(-\gamma_1)$, for any
 550 $j = 1, 2, \dots, l$,

$$\|D - D^*\|_{\max} \leq 2^j \bar{r} \Rightarrow \|D - D^*\|_{n_p, n_q} \leq 2^{j+1} \bar{r}.$$

551 Suppose that for some $j \leq l$, $\hat{D} \in B_{\max}(D^*, 2^j \bar{r}) \setminus B_{\max}(D^*, 2^{j-1} \bar{r})$, then by the results in Step 2,
 552 with probability at least $1 - 8 \exp(-\gamma_1)$,

$$\begin{aligned} c_0 \|\hat{D} - D^*\|_{\max}^2 &\leq (C_1 + C_2) 2^j \bar{r} \left(768 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{2\gamma_1}{n}} \right) \\ &\quad + C_0 \epsilon_n^2 + \sqrt{\frac{2\gamma_1}{n}} (C_1 + C_2) \epsilon_n + \frac{62(C_1 + C_2)M\gamma_1}{3n}. \end{aligned}$$

553 If

$$\frac{1}{c_0} (C_1 + C_2) \left(768 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log N}{N}} + \sqrt{2} \sqrt{\frac{\gamma_1}{N}} \right) \leq \frac{1}{8} 2^j \bar{r}, \quad (\text{A.15})$$

554 and

$$\frac{1}{c_0} \left[C_0 \epsilon_N^2 + \sqrt{2} (C_1 + C_2) \epsilon_N \sqrt{\frac{\gamma_1}{N}} + \frac{62(C_1 + C_2)M\gamma_1}{3N} \right] \leq \frac{1}{8} 2^{2j} \bar{r}^2, \quad (\text{A.16})$$

555 then

$$\|\hat{D} - D^*\|_{\max}^2 \leq 2^{2j-2} \bar{r}^2 \Leftrightarrow \|\hat{D} - D^*\|_{\max} \leq 2^{j-1} \bar{r}.$$

556 In short, to obtain this inequality, we need \bar{r} satisfying (A.15), (A.16) and $\bar{r} \geq$
 557 $\max(\sqrt{\log n/n}, 4\sqrt{3}M\sqrt{\gamma_1/n}/3, r_*)$. As $r_* \leq \kappa M \sqrt{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n/n}$, there exists a con-
 558 stant $C_* = C_*(c_0, C_0, C_1, C_2, M) = C'(\mu, \sigma, M)$ such that

$$\bar{r} = C_* \left(\sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{\gamma_1}{n}} + \epsilon_n \right)$$

559 satisfies all the requirements. As a result, with probability at least $1 - 10l \exp(-\gamma_1)$,

$$\|\hat{D} - D^*\|_{\max} \leq \bar{r} \text{ and } \|\hat{D} - D^*\|_{n_p, n_q} \leq 2\bar{r}.$$

560 Let $\gamma_1 = \log 10l + \gamma$, $l = \lfloor \log_2(2M/\sqrt{\log n/n}) \rfloor$, there exists $C = C(c_0, C_0, C_1, C_2, M) =$
 561 $C(\mu, \sigma, M)$ such that with probability at least $1 - \exp(-\gamma)$,

$$\|\hat{D} - D^*\|_{\max} \leq \bar{r} \leq C \left(\sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{\gamma}{n}} + \epsilon_n \right),$$

562 and

$$\|\hat{D} - D^*\|_{n_p, n_q} \leq 2C \left(\sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + \sqrt{\frac{\gamma}{n}} + \epsilon_n \right).$$

563 The proof of Theorem 1 is completed. \square

564 **Lemma A.3.** *The following excess risk decomposition always holds:*

$$\mathcal{B}_\psi(e^{\hat{D}}) - \mathcal{B}_\psi(e^{D^*}) = \left\{ \mathcal{B}_\psi(e^{\hat{D}}) - \inf_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{B}_\psi(e^D) \right\} + \left\{ \inf_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{B}_\psi(e^D) - \mathcal{B}_\psi(e^{D^*}) \right\} \quad (\text{A.17})$$

565 Under Assumptions 1 and 3, when $n \geq \text{Pdim}(\mathcal{F}_{\text{FNN}})$, there exist three constants C, C_0, C_* , with
 566 C, C_0 depending only on (μ, σ, M) and C_* depending only on (μ, σ) , such that

$$E_{p^*, q^*} \left\{ \mathcal{B}_\psi(e^{\hat{D}}) - \inf_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{B}_\psi(e^D) \right\} \leq C \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}}, \quad (\text{A.18})$$

567 and

$$E_{p^*, q^*} \|\hat{D} - D^*\|_p^2 \leq C_0 \sqrt{\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n}} + C_* e^{2M} \inf_{D \in \mathcal{F}_{\text{FNN}}} \|e^D - e^{D^*}\|_p^2.$$

568 *Proof of Lemma A.3.* To show (A.18) is the key step in the proof of this theorem, thus we focus on
 569 the proof of (A.18). Let

$$D_0 \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{B}_\psi(e^D).$$

570 Then,

$$\begin{aligned} & E_{p^*, q^*} \left\{ \mathcal{B}_\psi(e^{\hat{D}}) - \inf_{D \in \mathcal{F}_{\text{FNN}}} \mathcal{B}_\psi(e^D) \right\} \\ &= E_{p^*, q^*} \left\{ \mathcal{B}_\psi(e^{\hat{D}}) - \mathcal{B}_\psi(e^{D_0}) \right\} \\ &\leq E_{p^*, q^*} \left\{ \mathcal{B}_\psi(e^{\hat{D}}) - \hat{\mathcal{B}}_\psi(e^{\hat{D}}) + \hat{\mathcal{B}}_\psi(e^{\hat{D}}) - \hat{\mathcal{B}}_\psi(e^{D_0}) \right\} \\ &+ E_{p^*, q^*} \left\{ \hat{\mathcal{B}}_\psi(e^{D_0}) - \mathcal{B}_\psi(e^{D_0}) \right\} \\ &\leq 2E_{p^*, q^*} \left\{ \sup_{D \in \mathcal{F}_{\text{FNN}}} |\hat{\mathcal{B}}_\psi(e^D) - \mathcal{B}_\psi(e^D)| \right\}. \end{aligned} \quad (\text{A.19})$$

571 By the symmetrization technique, Talagrand's lemma, (A.12) and the fact that $\|D\|_\infty \leq M$ for any
 572 $D \in \mathcal{F}_{\text{FNN}}$, we can easily get the inequality (A.18) through (A.19). \square

573 *Proof of Theorem 2.* Theorem 2 is a direct corollary of Lemma A.3. We omit the details here. \square

574 *Proof of Theorem 3.* Since $D^* \in \mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$ and $a \in (0, 1]$,
 575 by Lemma 1, for the \mathcal{F}_{FNN} , a function class consists of ReLU FNN with width $\mathcal{W} = 38(\lfloor \beta \rfloor +$
 576 $1)^2 d^{\lfloor \beta \rfloor + 1} L \lceil \log_2(8L) \rceil$ and depth $\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 K \lceil \log_2(8K) \rceil$, where $K, L \in \mathbb{N}^+$, there exists
 577 a function $\phi_0 \in \mathcal{F}_{\text{FNN}}$ such that

$$\sup_{x \in [0, 1]^d \setminus H_{B, \delta}} |D^* - \phi_0| \leq 18M(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + (\beta \vee 1)/2} (KL)^{-\frac{2\beta}{d}}, \quad (\text{A.20})$$

578 where $H_{B, \delta} = \cup_{i=1}^d \{x = [x_1, \dots, x_d] : x_i \in \cup_{b=1}^{B-1} (b/B - \delta, b/B)\}$, $B = \lceil (KL)^{2/d} \rceil$, $\delta \in$
 579 $(0, 1/(3B)]$. As $D_{\text{NN}} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \|D - D^*\|_{\max}$, then

$$\|D_{\text{NN}} - D^*\|_{\max}^2 \leq \|\phi_0 - D^*\|_{\max}^2.$$

580 By the result in (A.20), for $I = p$ or q , we have

$$\begin{aligned} \|\phi_0 - D^*\|_I^2 &= \int_{[0, 1]^d \setminus H_{B, \delta}} |D^* - \phi_0|^2 I^*(x) dx + \int_{H_{B, \delta}} |D^* - \phi_0|^2 I^*(x) dx \\ &\leq 324M^2(\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)} (KL)^{-\frac{4\beta}{d}} + 4M^2 \int_{H_{B, \delta}} I^*(x) dx. \end{aligned}$$

581 As $p^*(\cdot), q^*(\cdot)$ are the density functions of some measures on $[0, 1]^d$ which are absolutely continuous
 582 with respect to the Lebesgue measure and δ can be arbitrarily small, $\int_{H_{B, \delta}} I_0(x) dx$ is also arbitrarily
 583 small. Thus we have

$$\|\phi_0 - D^*\|_I^2 \leq 324M^2(\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)} (KL)^{-\frac{4\beta}{d}}$$

584 and

$$\begin{aligned}
\|D_{\text{NN}} - D^*\|_{\max}^2 &\leq \|\phi_0 - D^*\|_{\max}^2 \\
&\leq 324M^2(\lfloor\beta\rfloor + 1)^4 d^{2\lfloor\beta\rfloor + (\beta\vee 1)} (KL)^{-\frac{4\beta}{d}} \\
&= 324M^2 C_1(\beta, d) (KL)^{-\frac{4\beta}{d}}.
\end{aligned}$$

585 By Corollary 1, there exists a constant C_1 only depending on (μ, σ, M) such that

$$\begin{aligned}
E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 &\leq C_1 \left(\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + \|D_{\text{NN}} - D^*\|_{\max}^2 \right) \\
&\leq C_1 \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + 324M^2 C_1(\beta, d) (KL)^{-\frac{4\beta}{d}} \right\} \\
&\leq 324M^2 C_1 \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + C_1(\beta, d) (KL)^{-\frac{4\beta}{d}} \right\}. \quad (\text{A.21})
\end{aligned}$$

586 This completes the proof of the first part of Theorem 3.

587 As for the second part of this theorem, based on Theorem 6 in Bartlett et al. (2019), for a specific
588 ReLU network f_ϕ , where ϕ contains the parameters in the network, there exists a universal constant
589 C_2 such that

$$\text{Pdim}(\mathcal{F}_{\text{FNN}}) \leq C_2 \mathcal{S} \mathcal{D} \log \mathcal{S},$$

590 where \mathcal{S} is the total number of parameters in the network f_ϕ . For a ReLU FNN with width
591 \mathcal{W} and depth \mathcal{D} , it can be easily checked that $\mathcal{S} = O(\mathcal{W}^2 \mathcal{D})$. Now for $\mathcal{W} = 114(\lfloor\beta\rfloor +$
592 $1)^2 d^{\lfloor\beta\rfloor + 1}$, $\mathcal{D} = 21(\lfloor\beta\rfloor + 1)^2 \left\lceil n^{\frac{d}{2(d+2\beta)}} \log_2 \left(8n^{\frac{d}{2(d+2\beta)}} \right) \right\rceil$, and \mathcal{W}, \mathcal{D} satisfy $O(\mathcal{W}^2 \mathcal{D}) =$
593 $O\left((\lfloor\beta\rfloor + 1)^6 d^{2\lfloor\beta\rfloor + 2} \left\lceil n^{\frac{d}{2(d+2\beta)}} \log^{-3} n \right\rceil \right)$, which means $L = 1, K = \left\lceil n^{\frac{d}{2(d+2\beta)}} \right\rceil$, and there exist
594 three universal constants C_3, C_4, C_5 such that

$$\begin{aligned}
&\frac{\mathcal{S} \mathcal{D} \log \mathcal{S} \log n}{n} \\
&\leq C_3 \left\{ (\lfloor\beta\rfloor + 1)^6 d^{2\lfloor\beta\rfloor + 2} \left\lceil n^{\frac{d}{2(d+2\beta)}} \log^{-3} n \right\rceil \right\} \\
&\quad \times \left(\log \left[C_3 \left\{ (\lfloor\beta\rfloor + 1)^6 d^{2\lfloor\beta\rfloor + 2} \left\lceil n^{\frac{d}{2(d+2\beta)}} \log^{-3} n \right\rceil \right\} \right] \right) \\
&\quad \times \left\{ 21(\lfloor\beta\rfloor + 1)^2 \left\lceil n^{\frac{d}{2(d+2\beta)}} \log_2 \left(8n^{\frac{d}{2(d+2\beta)}} \right) \right\rceil \log n / n \right\} \\
&\leq \frac{C_4}{n} \left\{ (\lfloor\beta\rfloor + 1)^8 d^{2\lfloor\beta\rfloor + 2} n^{\frac{2d}{2(d+2\beta)}} \log^{-1} n \right\} \\
&\quad \times \left\{ 6 \log(\lfloor\beta\rfloor + 1) + 2(\lfloor\beta\rfloor + 1) \log d + \frac{d}{2(d+2\beta)} \log n \right\} \\
&\leq \frac{C_4}{n} \left\{ (\lfloor\beta\rfloor + 1)^8 d^{2\lfloor\beta\rfloor + 2} n^{\frac{2d}{2(d+2\beta)}} \log^{-1} n \right\} \{6(\lfloor\beta\rfloor + 1) + 2(\lfloor\beta\rfloor + 1)d + \log n\} \\
&\leq C_5 (\lfloor\beta\rfloor + 1)^9 d^{2\lfloor\beta\rfloor + 3} n^{-\frac{2\beta}{d+2\beta}}. \quad (\text{A.22})
\end{aligned}$$

595 It follows from (A.21) that

$$\begin{aligned}
&E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \\
&\leq 324M^2 C_1 \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + C_1(\beta, d) (KL)^{-\frac{4\beta}{d}} \right\} \\
&\leq 324M^2 C_1 \left\{ \frac{C_2 \mathcal{S} \mathcal{D} \log \mathcal{S} \log n}{n} + C_1(\beta, d) (KL)^{-\frac{4\beta}{d}} \right\} \\
&\leq 324M^2 C_1 \left\{ C_2 C_5 (\lfloor\beta\rfloor + 1)^9 d^{2\lfloor\beta\rfloor + 3} n^{-\frac{2\beta}{d+2\beta}} + (\lfloor\beta\rfloor + 1)^4 d^{2\lfloor\beta\rfloor + (\beta\vee 1)} n^{-\frac{2\beta}{d+2\beta}} \right\} \\
&\leq 324M^2 C_1 (C_2 C_5 + 1) (\lfloor\beta\rfloor + 1)^9 d^{2\lfloor\beta\rfloor + (\beta\vee 3)} n^{-\frac{2\beta}{d+2\beta}}.
\end{aligned}$$

596 This completes the proof of the second part of Theorem 3. \square

597 *Proof of Theorem 4.* Based on Theorem 3.1 in Baraniuk & Wakin (2009), there exists a linear
 598 projection $A \in \mathbb{R}^{d_\delta \times d}$ such that $AA^T = dI_{d_\delta}/d_\delta$, where $I_{d_\delta} \in \mathbb{R}^{d_\delta \times d_\delta}$ is an identity matrix, and for
 599 any $x, y \in \mathcal{M}$,

$$(1 - \delta)\|x - y\|_2 \leq \|Ax - Ay\|_2 \leq (1 + \delta)\|x - y\|_2. \quad (\text{A.23})$$

600 Then we have

$$A(\mathcal{M}_\rho) \subseteq A([0, 1]^d) \subseteq \left[-\frac{d}{\sqrt{d_\delta}}, \frac{d}{\sqrt{d_\delta}}\right]^{d_\delta}.$$

601 Note that for any $z \in A(\mathcal{M})$, there exists a unique $x \in \mathcal{M}$ such that $z = Ax$. Otherwise, suppose we
 602 can find $x, x' \in \mathcal{M}, x \neq x'$ such that $z = Ax = Ax'$, then by (A.23), we know $\|x - x'\|_2 = 0$ and
 603 thus $x = x'$, which contradicts the assumption that $x \neq x'$. This uniqueness allows us to define a
 604 linear operator $\mathcal{S}\mathcal{L} : A(\mathcal{M}) \rightarrow \mathcal{M}$ such that $A[\mathcal{S}\mathcal{L}(z)] = z$. By (A.23), we have

$$(1 - \delta)\|\mathcal{S}\mathcal{L}(z_1) - \mathcal{S}\mathcal{L}(z_2)\|_2 \leq \|z_1 - z_2\|_2 \leq (1 + \delta)\|\mathcal{S}\mathcal{L}(z_1) - \mathcal{S}\mathcal{L}(z_2)\|_2.$$

605 This implies that the norm of $\mathcal{S}\mathcal{L}$ belongs to $[1/(1 + \delta), 1/(1 - \delta)]$. For the high-dimensional
 606 function $D^* : [0, 1]^d \rightarrow \mathbb{R}$ whose support is \mathcal{M}_ρ , it has an approximate low-dimensional representation
 607 \tilde{D}^* as follows:

$$\tilde{D}^*(z) = D^*(\mathcal{S}\mathcal{L}(z)), \quad \forall z \in A(\mathcal{M}).$$

608 As $D^* \in \mathcal{H}^\beta([0, 1]^d, M)$ with $\beta = k + a$ where $k \in \mathbb{N}^+$ and $a \in (0, 1]$, we have
 609 $\tilde{D}^* \in \mathcal{H}^\beta(A(\mathcal{M}), M/(1 - \delta)^\beta)$. By the extended version of Whitney's extension theorem
 610 in Fefferman (2006), since $A(\mathcal{M}) \subseteq A([0, 1]^d) \subseteq [-d/\sqrt{d_\delta}, d/\sqrt{d_\delta}]^{d_\delta}$, there exists $\tilde{D}_E^* \in$
 611 $\mathcal{H}^\beta\left([-d/\sqrt{d_\delta}, d/\sqrt{d_\delta}]^{d_\delta}, M/(1 - \delta)^\beta\right)$ such that $\tilde{D}_E^* \equiv \tilde{D}^*$ on $A(\mathcal{M})$. If $\mathcal{W} = 38(\lfloor \beta \rfloor +$
 612 $1)^2 d_\delta^{\lfloor \beta \rfloor + 1} L \lceil \log_2(8L) \rceil$ and $\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 K \lceil \log_2(8K) \rceil$, by the first result of Lemma 1,
 613 there exists a function ϕ_0 implemented by a ReLU network with width \mathcal{W} and depth \mathcal{D} such that

$$\sup_{z \in [0, 1]^{d_\delta} \setminus H_{B, \epsilon}^{d_\delta}} \left| \tilde{D}_E^* \left(\frac{2dz - d\mathbf{1}_{d_\delta}}{\sqrt{d_\delta}} \right) - \phi_0(z) \right| \leq \frac{18M}{(1 - \delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}}.$$

614 where $H_{B, \epsilon}^{d_\delta} = \cup_{i=1}^{d_\delta} \{x = [x_1, x_2, \dots, x_{d_\delta}] : x_i \in \cup_{b=1}^{B-1} (b/B - \epsilon, b/B)\}$ and $B =$
 615 $\lceil (KL)^{2/d} \rceil, \epsilon \in (0, 1/(3B))$. Thus

$$\begin{aligned} & \sup_{z \in \left[-\frac{d}{\sqrt{d_\delta}}, \frac{d}{\sqrt{d_\delta}}\right]^{d_\delta} \setminus \tilde{H}_{B, \epsilon}^{d_\delta}} \left| \tilde{D}_E^*(z) - \phi_0 \left(\frac{\sqrt{d_\delta}z + d\mathbf{1}_{d_\delta}}{2d} \right) \right| \\ & \leq \frac{18M}{(1 - \delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}}, \end{aligned}$$

616 where $\tilde{H}_{B, \epsilon}^{d_\delta} = \left\{ (2dt - d\mathbf{1}_{d_\delta})/\sqrt{d_\delta} : t \in H_{B, \epsilon}^{d_\delta} \right\}$.

617 Let $\tilde{\phi}_0(x) = \phi_0((\sqrt{d_\delta}Ax + d\mathbf{1}_{d_\delta})/(2d))$ and $\tilde{H}_{*B, \epsilon}^d =$
 618 $\left\{ x \in [0, 1]^{d \times d} : (\sqrt{d_\delta}Ax + d\mathbf{1}_{d_\delta})/(2d) \in H_{B, \epsilon}^{d_\delta} \right\}$. It can be easily checked that $\tilde{\phi}_0$ is also a function
 619 implemented by a ReLU network with the same structure as ϕ_0 , except that the input layer of $\tilde{\phi}_0$ has
 620 d units, instead of d_δ units. For any $x \in \mathcal{M}_\rho \setminus \tilde{H}_{*B, \epsilon}^d, Ax \in [-d/\sqrt{d_\delta}, d/\sqrt{d_\delta}]^{d_\delta} \setminus \tilde{H}_{B, \epsilon}^{d_\delta}$ and there
 621 exists a $x' \in \mathcal{M}$ satisfying $\|x - x'\|_2 \leq \rho$. Since $\tilde{D}_E^* \in \mathcal{H}^\beta\left([-d/\sqrt{d_\delta}, d/\sqrt{d_\delta}]^{d_\delta}, M/(1 - \delta)^\beta\right)$

622 and $D^* \in \mathcal{H}^\beta([0, 1]^d, M)$,

$$\begin{aligned}
& |\tilde{\phi}_0(x) - D^*(x)| \\
& \leq |\tilde{\phi}_0(x) - \tilde{D}_E^*(Ax)| + |\tilde{D}_E^*(Ax) - \tilde{D}_E^*(Ax')| + |\tilde{D}_E^*(Ax') - D^*(x)| \\
& \leq \frac{18M}{(1-\delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}} + \frac{M}{(1-\delta)^\beta} \|Ax' - Ax\|_2 + \rho M \\
& \leq \frac{18M}{(1-\delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}} + \frac{M\sqrt{d}}{(1-\delta)^\beta \sqrt{d_\delta}} \rho + \rho M \\
& \leq \frac{18M}{(1-\delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}} + \frac{2M\sqrt{d}}{(1-\delta)^\beta \sqrt{d_\delta}} \rho \\
& \leq \frac{20M}{(1-\delta)^\beta} (\lfloor \beta \rfloor + 1)^2 (2d)^\beta d_\delta^{\lfloor \beta \rfloor + (\beta \vee 1 + \beta)/2} (KL)^{-\frac{2\beta}{d_\delta}}, \tag{A.24}
\end{aligned}$$

623 where the last inequality holds when $\rho \leq (\lfloor \beta \rfloor + 1)^2 2^\beta d^{\beta - \frac{1}{2}} d_\delta^{\lfloor \beta \rfloor + (\beta - 1/2) \vee (1/2)} (KL)^{-\frac{2\beta}{d_\delta}}$. As
624 $D_{\text{NN}} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \|D - D^*\|_{\max}$,

$$\|D_{\text{NN}} - D^*\|_{\max}^2 \leq \|\tilde{\phi}_0 - D^*\|_{\max}^2.$$

625 By the result in (A.24), for $I = p$ or q , it holds

$$\begin{aligned}
\|\tilde{\phi}_0 - D^*\|_I^2 &= \int_{[0,1]^d \setminus H_{B,\delta}} |D^* - \tilde{\phi}_0|^2 I^*(x) dx + \int_{H_{B,\delta}} |D^* - \tilde{\phi}_0|^2 I^*(x) dx \\
&\leq \frac{400M^2}{(1-\delta)^{2\beta}} (\lfloor \beta \rfloor + 1)^4 (2d)^{2\beta} d_\delta^{\beta \vee 1 + 3\beta} (KL)^{-\frac{4\beta}{d_\delta}} + \frac{4M^2}{(1-\delta)^{2\beta}} \int_{\tilde{H}_{*B,\epsilon}^d} I^*(x) dx.
\end{aligned}$$

626 As $p^*(\cdot), q^*(\cdot)$ are the density functions of some measures on $[0, 1]^d$ which are absolutely continuous
627 w.r.t the Lebesgue measure and ϵ can be arbitrarily small for the given δ , $\int_{\tilde{H}_{*B,\epsilon}^d} I_0(x) dx$ is also
628 arbitrarily small for the given δ . Thus we have

$$\|\tilde{\phi}_0 - D^*\|_I^2 \leq \frac{400M^2}{(1-\delta)^{2\beta}} (\lfloor \beta \rfloor + 1)^4 (2d)^{2\beta} d_\delta^{\beta \vee 1 + 3\beta} (KL)^{-\frac{4\beta}{d_\delta}}$$

629 and

$$\begin{aligned}
\|D_{\text{NN}} - D^*\|_{\max}^2 &\leq \|\phi_0 - D^*\|_{\max}^2 \\
&\leq \frac{400M^2}{(1-\delta)^{2\beta}} (\lfloor \beta \rfloor + 1)^4 (2d)^{2\beta} d_\delta^{\beta \vee 1 + 3\beta} (KL)^{-\frac{4\beta}{d_\delta}} \\
&= \frac{400M^2}{(1-\delta)^{2\beta}} C_2(\beta, d, d_\delta) (KL)^{-\frac{4\beta}{d_\delta}}.
\end{aligned}$$

630 By Corollary 1, there exists a constant C_1 only depending on (μ, σ, M) , such that

$$\begin{aligned}
& E_{p^*, q^*} \|\hat{D} - D^*\|_{\max}^2 \\
& \leq C_1 \left(\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + \|D_{\text{NN}} - D^*\|_{\max}^2 \right) \\
& \leq C_1 \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + \frac{400M^2}{(1-\delta)^{2\beta}} C_2(\beta, d, d_\delta) (KL)^{-\frac{4\beta}{d_\delta}} \right\} \\
& \leq \frac{400M^2 C_1}{(1-\delta)^{2\beta}} \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + C_2(\beta, d, d_\delta) (KL)^{-\frac{4\beta}{d_\delta}} \right\}. \tag{A.25}
\end{aligned}$$

631 For

$$\mathcal{W} = 114(\lfloor \beta \rfloor + 1)^2 d_\delta^{\lfloor \beta \rfloor + 1},$$

632

$$\mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 \left\lceil n^{\frac{d_\delta}{2(d_\delta + 2\beta)}} \log_2 \left(8n^{\frac{d_\delta}{2(d_\delta + 2\beta)}} \right) \right\rceil,$$

633 and \mathcal{W}, \mathcal{D} satisfy

$$\mathcal{O}(\mathcal{W}^2 \mathcal{D}) = \mathcal{O} \left(([\beta] + 1)^6 d_\delta^{2[\beta]+2} \left[n^{\frac{d_\delta}{2(d_\delta+2\beta)}} \log^{-3} n \right] \right),$$

634 along the derivation of (A.22), there exists a universal constants C^* such that

$$\frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} \leq C^* ([\beta] + 1)^9 d_\delta^{2[\beta]+3} n^{-\frac{2\beta}{d_\delta+2\beta}}.$$

635 Based on the result of (A.25),

$$\begin{aligned} & E_{p^*, q^*} \|\widehat{D} - D^*\|_{\max}^2 \\ & \leq \frac{400M^2 C_1}{(1-\delta)^{2\beta}} \left\{ \frac{\text{Pdim}(\mathcal{F}_{\text{FNN}}) \log n}{n} + C_2(\beta, d, d_\delta) (KL)^{-\frac{4\beta}{d_\delta}} \right\} \\ & \leq \frac{400M^2 C_1}{(1-\delta)^{2\beta}} \left\{ C^* ([\beta] + 1)^9 d_\delta^{2[\beta]+3} n^{-\frac{2\beta}{d_\delta+2\beta}} + C_2(\beta, d, d_\delta) n^{-\frac{2\beta}{d_\delta+2\beta}} \right\} \\ & \leq \frac{800M^2 C_1 C^*}{(1-\delta)^{2\beta}} ([\beta] + 1)^9 \max \left\{ d_\delta^{2[\beta]+3}, (2d)^{2\beta} d_\delta^{\beta \vee 1 + 3\beta} \right\} n^{-\frac{2\beta}{d_\delta+2\beta}} \\ & = \frac{800M^2 C_1 C^* C_3(\beta, d, d_\delta)}{(1-\delta)^{2\beta}} n^{-\frac{2\beta}{d_\delta+2\beta}}. \end{aligned}$$

636 This completes the proof of the theorem and (12). \square

637 *Proof of Proposition 1.* For $k = 0, \dots, K-2$, the densities $q_k(\cdot), q_{k+1}(\cdot)$ of the synthetic data
638 $\{Z_{k,j}\}_{j=1}^n$ and $\{Z_{k+1,j}\}_{j=1}^n$ satisfy

$$\frac{q_k(t)}{q_{k+1}(t)} = \frac{(1-\alpha_k)q^*(z) + \alpha_k p^*(z)}{(1-\alpha_{k+1})q^*(z) + \alpha_{k+1} p^*(z)} \in \left[\frac{(1-e^{-M})\alpha_k + e^{-M}}{(1-e^{-M})\alpha_{k+1} + e^{-M}}, \frac{1-\alpha_k}{1-\alpha_{k+1}} \right].$$

639 As $\|f\|_2 = (\int_{\mathcal{Z}} f^2(x) dx)^{\frac{1}{2}}$, then for any density g satisfying $g \geq c$, $\|f\|_2 = (\int_{\mathcal{Z}} f^2(x) dx)^{\frac{1}{2}} \leq$
640 $(\int_{\mathcal{Z}} f^2(x) g(x) / c dx)^{\frac{1}{2}} = \|f\|_g / \sqrt{c}$. Using an appropriate $\mathcal{F}_{\text{FNN}}^0$ whose element D satisfies $\|D\|_\infty \leq$
641 M_0 , for the direct estimate \widehat{D}_{SRE} , as $\log(q^*/p^*)$ is only bounded from below by $-M_0$, by Theorem
642 2, we have

$$\limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{\text{SRE}} - D^*\|_2 \leq e^{M_0} C_*(\mu, \sigma, c_1) \|R^* - R_{M_0}^*\|_p.$$

643 For $k = 0, 1, \dots, K-2$, as $|\log\{q_k(t)/q_{k+1}(t)\}|$ is bounded by M_0 , by Corollary 1, we have

$$\limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_k - D_k^*\|_2 = 0.$$

644 Let $R_{K-1, M_0}^* = (1-\alpha_{K-1})R_{M_0}^* + \alpha_{K-1}$. As the logarithm of $R_{K-1}^* = (1-\alpha_{K-1})q^*/p^* + \alpha_{K-1}$
645 is also only bounded from below by $-M_0$, again, by Theorem 2,

$$\begin{aligned} \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{K-1} - D_{K-1}^*\|_2 & \leq e^{M_0} C_*(\mu, \sigma, c_1) \|R_{K-1}^* - R_{K-1, M_0}^*\|_p \\ & = (1-\alpha_{K-1}) e^{M_0} C_*(\mu, \sigma, c_1) \|R^* - R_{M_0}^*\|_p. \end{aligned}$$

646 Thus

$$\begin{aligned} \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{\text{TRE}} - D^*\|_2 & \leq \sum_{k=0}^{K-1} \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_k - D_k^*\|_2 \\ & = \limsup_{n \rightarrow \infty} E_{p^*, q^*} \|\widehat{D}_{K-1} - D_{K-1}^*\|_2 \\ & \leq (1-\alpha_{K-1}) e^{M_0} C_*(\mu, \sigma, c_1) \|R^* - R_{M_0}^*\|_p. \end{aligned}$$

647 This completes the proof of Proposition 1. \square

648 **A.2 Examples of Hölder function class**

649 Let p^* be the density function of a truncated d -dimensional multivariate Gaussian with mean zero
650 and covariance matrix $\Sigma_p \in \mathbb{R}^{d \times d}$ in $[0, 1]^d$. That means

$$p^*(z) = \exp(-z' \Sigma_p^{-1} z / 2) / c(\Sigma_p), \quad c(\Sigma_p) = \int_{[0,1]^d} \exp(-t' \Sigma_p^{-1} t / 2) dt, \quad z \in [0, 1]^d.$$

651 Similarly, let

$$q^*(z) = \exp(-z' \Sigma_q^{-1} z / 2) / c(\Sigma_q)$$

652 for some positive definite matrix Σ_q . For any matrix $A \in \mathbb{R}^{d \times d}$, $A_{i\cdot}$ is the i th row of A for
653 $i = 1, 2, \dots, d$ and

$$\|A\|_{2,\infty} := \sup_{\|z\|_\infty \leq 1} \|Az\|_2.$$

654 Then,

$$D^*(z) = \log \frac{q^*(z)}{p^*(z)} = \frac{1}{2} z' (\Sigma_p^{-1} - \Sigma_q^{-1}) z + \log(c(\Sigma_p) - c(\Sigma_q)), \quad z \in [0, 1]^d.$$

655 Let $M = \max \left\{ \frac{1}{2} (\|\Sigma_p^{-1/2}\|_{2,\infty}^2 + \|\Sigma_q^{-1/2}\|_{2,\infty}^2) + |\log[c(\Sigma_p) - c(\Sigma_q)]|, \|(\Sigma_p^{-1} - \Sigma_q^{-1})_{i\cdot}\|_2, i = 1, 2, \dots, d \right\}$.

656 It is straightforward to check that

$$D^* \in \mathcal{H}^2([0, 1]^d, M).$$

657 It implies the Hölder smoothness parameter β is 2 for this example.

658 Moreover, the truncated multivariate Gaussian distributions considered above are special cases of the
659 exponential distribution class defined below. Define the density function class

$$\text{Exp}(\beta, B) := \left\{ p(z) = \exp(g(z)) / c_g : z \in [0, 1]^d, c_g = \int_{[0,1]^d} \exp(g(t)) dt, g \in \mathcal{H}^\beta([0, 1]^d, B) \right\}.$$

660 Suppose that $\Sigma \in \mathbb{R}^{d \times d}$ is positive definite and let $g(z) = z' \Sigma z / 2$. Then, $g \in \mathcal{H}^2([0, 1]^d, M_\Sigma)$,
661 where $M_\Sigma = \max \left\{ \frac{1}{2} (\|\Sigma^{1/2}\|_{2,\infty}^2, \|\Sigma_{i\cdot}\|_2, i = 1, 2, \dots, d) \right\}$. If $p^*, q^* \in \text{Exp}(\beta, B)$, we have
662 $D^*(z) = \log[q^*(z)/p^*(z)] \in \mathcal{H}^\beta([0, 1]^d, 4B)$.

663 **A.3 Extension to unbounded support case**

664 In fact, our Theorem 1, Corollary 1 and Theorem 2 do not rely on the hypercube assumption. To
665 relax the hypercube assumption to allow unbounded support, we need to study the upper bound
666 for the approximation error $\|D_{\text{NN}} - D^*\|_{\max}$ carefully. With unbounded support, we may bound
667 $\|D_{\text{NN}} - D^*\|_{\max}$ by the truncation technique under some additional assumptions, at a small price of
668 an additional logarithm term in the error bound.

669 Specifically, when the pdfs are supported on \mathbb{R}^d , to bound the approximation error as in Theorem
670 3, aside from Assumptions 1-2 and the Hölder class assumption, we need to further assume that
671 $\max\{E_{p^*} I(\|Z\|_\infty \geq \log n), E_{q^*} I(\|Z\|_\infty \geq \log n)\} \leq n^{-\frac{2\beta}{d+2\beta}}$. For $I = p$ or q , and any $D \in \mathcal{F}_{\text{FNN}}$,
672 where \mathcal{F}_{FNN} is the function class of ReLU FNNs with width \mathcal{W} and depth \mathcal{D} specified by

$$\mathcal{W} = 114(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1}, \quad \mathcal{D} = 21(\lfloor \beta \rfloor + 1)^2 \left\lceil n^{\frac{d}{2(d+2\beta)}} \log_2 \left(8n^{\frac{d}{2(d+2\beta)}} \right) \right\rceil,$$

673 we have

$$\begin{aligned} & E_{I^*} [D(Z) - D^*(Z)]^2 \\ & \leq E_{I^*} [\{D(Z) - D^*(Z)\}^2 I(\|Z\|_\infty \geq \log n)] + E_{I^*} [\{D(Z) - D^*(Z)\}^2 I(\|Z\|_\infty \leq \log n)] \\ & \leq 4M^2 E_{I^*} I(\|Z\|_\infty \geq \log n) + E_{I^*} [\{D(Z) - D^*(Z)\}^2 I(\|Z\|_\infty \leq \log n)], \end{aligned}$$

674 where the second inequality follows from the facts that $\|D^*\|_\infty \leq M, \|D\|_\infty \leq M$ under Assumption
675 2. Since $D^* \in \mathcal{H}^\beta(\mathbb{R}^d, M)$, $D^*(2t \log n - \log n \mathbf{1}_d) \in \mathcal{H}^\beta([0, 1]^d, (2 \log n)^{\lfloor \beta \rfloor} M)$ as a function of
676 t , where $\mathbf{1}_d$ is the d -dimensional all-one vector. By Lemma 1, there exists a function $\phi_0 \in \mathcal{F}_{\text{FNN}}$ such
677 that

$$\sup_{t \in [0, 1]^d \setminus H_{B, \delta}} |D^*(2t \log n - \log n \mathbf{1}_d) - \phi_0| \leq 18(2 \log n)^{\lfloor \beta \rfloor} M (\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + 1} n^{-\frac{\beta}{d+2\beta}},$$

678 where $H_{B,\delta} = \cup_{i=1}^d \{t = [t_1, \dots, t_d] : t_i \in \cup_{b=1}^{B-1} (b/B - \delta, b/B)\}$, $B = \lceil n^{\frac{1}{d+2\beta}} \rceil$, $\delta \in$
679 $(0, 1/(3B)]$. Thus

$$\sup_{z \in [-\log n, \log n]^d \setminus \tilde{H}_{B,\epsilon}^d} \left| D^*(z) - \phi_0 \left(\frac{z + \log n \mathbf{1}_d}{2 \log n} \right) \right| \leq 18(2 \log n)^{\lfloor \beta \rfloor} M(\lfloor \beta \rfloor + 1)^2 d^{\lfloor \beta \rfloor + (\beta \vee 1)/2} n^{-\frac{\beta}{d+2\beta}},$$

680 where $\tilde{H}_{B,\delta}^d = \{2t \log n - \log n : t \in H_{B,\delta}^d\}$. Let $\tilde{\phi}_0(z) = \phi_0 \left(\frac{z + \log n \mathbf{1}_d}{2 \log n} \right) \in \mathcal{F}_{\text{FNN}}$. As δ can be
681 arbitrarily small, it then follows from similar lines as in the proof of Theorem 3 that

$$E_{I^*} [\{\tilde{\phi}_0(Z) - D^*(Z)\}^2 I(\|Z\|_\infty \leq \log n)] \leq 324M^2(\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)} (2 \log n)^{2\lfloor \beta \rfloor} n^{-\frac{2\beta}{d+2\beta}}.$$

682 Since $D_{\text{NN}} \in \arg \min_{D \in \mathcal{F}_{\text{FNN}}} \|D - D^*\|_{\max}$, we have

$$\begin{aligned} \|D_{\text{NN}} - D^*\|_{\max}^2 &\leq \|\tilde{\phi}_0 - D^*\|_{\max}^2 \\ &\leq \max_{I=p,q} \{4M^2 E_{I^*} I(\|Z\|_\infty \geq \log n) + E_{I^*} [\{\tilde{\phi}_0(Z) - D^*(Z)\}^2 I(\|Z\|_\infty \leq \log n)]\} \\ &\leq 328M^2(\lfloor \beta \rfloor + 1)^4 d^{2\lfloor \beta \rfloor + (\beta \vee 1)} (2 \log n)^{2\lfloor \beta \rfloor} n^{-\frac{2\beta}{d+2\beta}}. \end{aligned}$$

683 Compared with the upper bound of the approximation error in Theorem 3, when the pdfs are supported
684 on \mathbb{R}^d (unbounded case), a similar approximation error upper bound can be derived with an additional
685 logarithmic factor $(2 \log n)^{2\lfloor \beta \rfloor}$.

686 A.4 Simulation setting and implementation details

687 Our simulation settings are as follows.

- 688 • Beta setting: Let $Z = (Z_1, Z_2, \dots, Z_p)^\top \in \mathbb{R}^p$ be a random vector of interest, where
689 Z_1, Z_2, \dots, Z_p are i.i.d. random variables following Beta distribution, denoted by
690 $\text{Beta}(\alpha, \beta)$. Set p^* as the p.d.f of $\text{Beta}(2.2, 1.5)$ and q^* as the p.d.f of $\text{Beta}(2, 2)$. In
691 this setting, we set $p = 5$.
- 692 • Normal setting: Let $Z = (Z_1, Z_2, \dots, Z_d, Z_{d+1}, Z_{d+2}, \dots, Z_{2d})^\top \in \mathbb{R}^{2d}$ be some random
693 vector of interest. Let p^* be the p.d.f of $N(0, I_{2d})$ and q^* be the p.d.f of $N(0, \Sigma(\rho))$, where
694 $\Sigma(\rho) = (\sigma_{i,j}^\rho) \in \mathbb{R}^{2d \times 2d}$ and

$$\sigma_{i,j}^\rho = \begin{cases} 1, & i = j; \\ \rho, & |i - j| = d, i, j = 1, 2, \dots, 2d; \\ 0, & \text{otherwise.} \end{cases}$$

695 In this setting, we set $d = 5$ and $\rho = 0.9$.

696 We apply the Adam algorithm (Kingma & Ba, 2014) in Pytorch with a learning rate $lr = 0.0001$ and
697 a weight decay parameter $wd = 0.0001$. A neural network with 2 hidden layers with widths (64, 64)
698 and ReLU activation function, is used in the experiment. The maximum number of epoches is 20000.
699 In this experiment, the training data size n is 5000 (10000). A validation data is used. The batch size
700 is 500 (1000), and an early-stopping technique is applied with $patience = 100$ for Beta setting and
701 $patience = 1000$ for Normal setting, where $patience$ is the number of epochs until termination if no
702 improvement is made on the validation dataset. The experiment is conducted on a laptop with an
703 *Intel(R) Core(TM) i7-8750H @ 2.20GHz* CPU having 6 cores. We use the LR-Bregman divergence in
704 this example. For the sequence $0 = \alpha_0 < \alpha_1 < \dots < \alpha_{K-1} < \alpha_K = 1$, we use the linearly spaced
705 α_k 's, that is $\alpha_k = k/K$, $k = 0, 1, 2, \dots, K$.

706 A.5 The MNIST dataset

707 We now apply the proposed mixing chain (13) for density ratio estimation to the MNIST dataset
708 (LeCun et al., 2010). In the implementation, to accelerate the computation, we use the subsampling
709 method with a training subsample size of 20,000 and a relatively small DenseNet network structure
710 (Huang et al., 2017); see Table A.1 for the specification of the network architectures. Similarly
711 to the results in Table 1 of Rhodes et al. (2020), we calculate the average negative log-likelihood
712 (ANLL) in bits per dimension (bpd, smaller is better). We denote the estimate based on the proposed

713 mixing chain (13) with the chain length B by "mTRE- B ". The batch size is 512, $lr = 0.001$ and
 714 $wd = 0.0001$. The maximum number of epoches is 1000. The reference distribution for our mTRE
 715 is taken to be the standard Gaussian distribution. Here, the reference distribution is the same as the
 716 noise distribution in the MNIST experiments of (Rhodes et al., 2020). We obtain the averaged ANLLs
 717 and their empirical standard errors for mTRE-5 and mTRE-10 over 5 random training subsamples.
 718 As a comparison, we use the results with the Gaussian noise for the direct single ratio estimate and
 719 the direct estimate based on the original convolution chain (cTRE) obtained from Table 1 in (Rhodes
 720 et al., 2020) as the benchmarks. Note that cTRE and the direct single ratio estimate are based on the
 721 full training sample, where the sample size is 60,000. The result for the cTRE presented in Table A.2
 722 is the best one among the cTRE's with the chain length $B \in \{5, 10, 15, 20, 25, 30\}$ in Table 1 in the
 723 online supplemental of Rhodes et al. (2020). Our results are presented in Table A.2.

724 From Table A.2, we see that mTRE is significantly better than the single ratio estimate and comparable
 725 with cTRE. The difference between the results from mTRE and cTRE is not statistically significant.
 726 We note that the training sample size for mTRE we used is restricted to 20,000, due to the memory
 727 limitation of the laptop used in the computation. In comparison, the sample size for cTRE is 60,000.

Table A.1: Architecture for mTRE

Layers	Details	Output size
Convolution	3×3 Conv	$12 \times 28 \times 28$
Transition Layer 1	BN, ReLU, 2×2 Average Pool, 1×1 Conv	$12 \times 14 \times 14$
Dense Block 1	BN, 1×1 Conv, BN, 3×3 Conv	$24 \times 14 \times 14$
Transition Layer 1	BN, ReLU, 2×2 Average Pool, 1×1 Conv	$12 \times 7 \times 7$
Dense Block 1	BN, 1×1 Conv, BN, 3×3 Conv	$24 \times 7 \times 7$
Pooling	BN, ReLU, 7×7 Average Pool, Reshape	24
Fully connected	Linear	1

Table A.2: Average negative log-likelihood (ANLL) in bits per dimension (bpd, smaller is better). For the proposed mixing chain estimate with the chain length B (mTRE- B), the ANLLs are averaged over 5 random training subsamples, where the subsample size is 20,000, and the corresponding standard errors are in parentheses. The cTRE is the direct estimate based on the original convolution chain (cTRE). The results for the direct single ratio estimate and the direct cTRE are obtained from Table 1 in the seminal paper (Rhodes et al., 2020) and we use them as the benchmarks. The cTRE and the direct single ratio estimate are based on the full training sample, where the sample size is 60,000.

Estimator	mTRE-5	mTRE-10	Direct Single ratio	Direct cTRE
ANLL	1.40 (0.0045)	1.39 (0.0077)	1.96	1.39