

# 000 001 BOUNDS ON $L_p$ ERRORS IN DENSITY RATIO ESTIMA- 002 TION VIA $f$ -DIVERGENCE LOSS FUNCTIONS 003 004

005 **Anonymous authors**

006 Paper under double-blind review

## 007 008 ABSTRACT 009

011 Density ratio estimation (DRE) is a fundamental machine learning technique for  
012 identifying relationships between two probability distributions.  $f$ -divergence loss  
013 functions, derived from variational representations of  $f$ -divergence, are com-  
014 monly employed in DRE to achieve state-of-the-art results. This study presents  
015 a novel perspective on DRE using  $f$ -divergence loss functions by deriving the up-  
016 per and lower bounds on  $L_p$  errors. These bounds apply to any estimator within a  
017 class of Lipschitz continuous estimators, irrespective of the specific  $f$ -divergence  
018 loss functions utilized. The bounds are formulated as a product of terms that in-  
019 clude the data dimension and the expected value of the density ratio raised to the  
020 power of  $p$ . Notably, the lower bound incorporates an exponential term depen-  
021 dent on the Kullback–Leibler divergence, indicating that the  $L_p$  error significantly  
022 increases with the Kullback–Leibler divergence for  $p > 1$ , and this increase be-  
023 comes more pronounced as  $p$  increases. Furthermore, these theoretical findings  
024 are substantiated through numerical experiments.

## 025 1 INTRODUCTION

028 Density ratio estimation (DRE) is a key technique in machine learning that calculates the density  
029 ratio  $r^*(\mathbf{x}) = q(\mathbf{x})/p(\mathbf{x})$  between two probability distributions based on samples drawn indepen-  
030 dently from  $p$  and  $q$ . DRE is integral to various machine learning methods such as generative mod-  
031 eling (Goodfellow et al., 2014; Nowozin et al., 2016; Uehara et al., 2016), mutual information esti-  
032 mation and representation learning (Belghazi et al., 2018; Hjelm et al., 2018), energy-based model-  
033 ing (Gutmann & Hyvärinen, 2010), and covariate shift and domain adaptation (Shimodaira, 2000;  
034 Huang et al., 2006).

035 Recent advancements in DRE have been driven by neural network-based methods, which utilize  
036 neural networks as density ratio estimators. These methods employ loss functions derived from  
037 variational representations of  $f$ -divergence (Nguyen et al., 2010; Sugiyama et al., 2012), where the  
038 optimal function corresponds to the density ratio through the Legendre transform, achieving state-  
039 of-the-art results.

040 Amidst their success, ongoing research has started to elucidate the theoretical relationship between  
041 the optimization of  $f$ -divergence loss functions and DRE accuracy. For integral probability metric  
042 (IPM) loss functions, the upper and lower bounds of the  $L_p$  error in DRE have been established as  
043 the minimax bounds of their optimization(Liang, 2017; Niles-Weed & Berthet, 2022). More recent  
044 studies have focused on  $f$ -divergence loss functions to derive the upper bounds (Belomestny et al.,  
045 2021) and the minimax upper and lower bounds for the optimization of Shannon divergence loss  
046 (Belomestny et al., 2021; Puchkin et al., 2024).

047 However, several aspects of this relationship remain unresolved. First, the minimax lower bounds  
048 do not represent the true lower bound of estimation accuracy for the actual density ratio. Second,  
049 the connection between the true magnitudes of  $f$ -divergences and the sample size requirements  
050 for DRE using divergence loss functions is not completely understood. Specifically, the impact of  
051 the true amount of Kullback–Leibler (KL) divergence on the sample size needed for DRE using  
052 the KL-loss function is unclear, despite known exponential increases in sample size requirements  
053 for KL-divergence estimation as the true KL-divergence widens (Poole et al., 2019; Song & Ermon,  
2019; McAllester & Stratos, 2020). Finally, it is not understood whether the  $L_p$  errors, e.g., the

054 root mean square errors (RMSE), of DRE are statistically equivalent, regardless of the choice of  
 055  $f$ -divergence loss function, such as the total variation loss or the KL-divergence loss function.  
 056

057 This study aims to address uncertainties in DRE using  $f$ -divergence loss functions by deriving the  
 058 upper and lower bounds that are independent of the choice of  $f$ -divergence. However, the theoretical  
 059 optimization of  $f$ -divergence loss functions is challenging owing to their reliance on sample sets  
 060 from two distributions. The lack of overlap in these sample sets leads to unstable optimization  
 061 points, causing the losses to fall below their theoretical optimal values. Practically, this issue is  
 062 often mitigated by implementing early stopping while monitoring validation losses.

063 We integrate this practical approach into our theoretical analysis framework through a conceptual  
 064 reformulation of the loss functions, thus bridging the gap between practical and theoretical behav-  
 065 iors of these functions. Subsequently, we derive upper and lower bounds for the  $L_p$  error in DRE  
 066 by optimizing  $f$ -divergence loss functions. These bounds are derived from the expectation of the  
 067 distance between the nearest neighbors in observations, assuming the  $L$ -Lipschitz continuity of the  
 068 energy function of the distributions and the compactness of the support.

069 The upper and lower bounds are formulated as a product of terms involving the data dimension and  
 070 the expectation of the density ratio raised to the power of  $p$ . Notably, the lower bound includes  
 071 an exponential term of the KL-divergence, indicating that the  $L_p$  error significantly increases as  
 072 the KL-divergence increases for  $p > 1$ , with the rate of increase accelerating for larger values of  
 073  $p$ . These bounds are applicable to a group of Lipschitz continuous estimators, irrespective of the  
 074 specific  $f$ -divergence loss functions employed. The theoretical implications are validated through  
 075 numerical experiments.

076 To summarize, the key contributions of this study are as follows: (1) We provide common upper  
 077 and lower bounds for the  $L_p$  error in DRE through optimizations of variational representations of  
 078  $f$ -divergences, introducing a novel understanding of DRE using  $f$ -divergence loss functions. (2) We  
 079 empirically investigate the relationship between KL-divergence, data dimension, and the estimation  
 080 accuracy of DRE through optimizations of variational representations of  $f$ -divergences. Specifically,  
 081 we discover that the  $L_p$  error significantly increases with the rise in KL-divergence when  $p > 1$ , and  
 082 this increase is exacerbated by the magnitude of the order  $p$ .

083 **Related Work.** This study provides upper and lower bounds on convergence rates for nonparamet-  
 084 ric density ratio estimation using  $f$ -divergence optimization. Relevant prior work includes studies  
 085 on the minimax convergence rates for density estimation within the context of GAN optimization,  
 086 specifically for Wasserstein GANs (Arjovsky & Bottou, 2017) and vanilla GANs (Goodfellow et al.,  
 087 2014). For Wasserstein GAN optimization, Liang (2017) and Singh & Póczos (2018) established  
 088 the minimax convergence rates for the IPW loss, which encompasses the total variation among  $f$ -  
 089 divergences. Additionally, Niles-Weed & Berthet (2022) extended these results to the Wasserstein- $p$   
 090 distance for  $p > 1$ . In the context of vanilla GAN optimization, Belomestny et al. (2021) and  
 091 Puchkin et al. (2024) presented minimax upper and lower convergence rates for the Shannon  
 092 divergence loss, providing an upper bound for the  $L_2$  error. Beyond GAN-related research, Nguyen et al.  
 093 (2010) presented an upper bound for the Hellinger distance in DRE using the KL-divergence loss,  
 094 thereby providing a minimax upper bound for the  $L_1$  error in DRE. Additionally, foundational work  
 095 by Stone (1980) established a minimax convergence rate for nonparametric regression, which is also  
 096 applicable to an upper bound for the  $L_1$  error in nonparametric density estimation.

## 097 2 PRELIMINARIES: NOTATION, SETUP, AND $f$ -DIVERGENCE LOSS 098 FUNCTIONS

101 In this section, we introduce the notation, problem setup, and the variational representation of  $f$ -  
 102 divergence, along with the corresponding loss functions that underpin the analysis in subsequent  
 103 sections.

### 105 2.1 NOTATION, PRELIMINARY CONCEPTS, AND SETUP

106 **Notation.** Random variables are denoted by uppercase letters, such as  $X$ . Lowercase letters repre-  
 107 sent specific values of these random variables; for instance,  $x$  denotes a value of the random variable

X. Boldface letters,  $\mathbf{X}$  and  $\mathbf{x}$ , denote the sets of random variables and their corresponding values, respectively.  $\|\mathbf{y} - \mathbf{x}\|_\infty$  denotes the maximum norm in  $\mathbb{R}^d$ . i.e.,  $\|\mathbf{y} - \mathbf{x}\|_\infty = \max_{1 \leq i \leq d} |y_i - x_i|$  for  $\mathbf{y} = (y_1, y_2, \dots, y_d)$  and  $\mathbf{x} = (x_1, x_2, \dots, x_d)$ .  $\text{diag}(\Omega)$  denotes the diameter of  $\Omega$ . Specifically, let  $\text{diag}(\mathcal{B}) = \inf_{r \in \mathbb{R}} \{\mathcal{B} \subseteq \Delta(\mathbf{a}, r) \mid \exists \mathbf{a} \in \mathcal{B}\}$ , where  $\Delta(\mathbf{a}, r)$  denotes the  $d$ -dimensional interval centered at  $\mathbf{a}$  with each side of length  $r$ :  $\Delta(\mathbf{a}, r) = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \mathbf{a}\|_\infty < r/2\}$ .  $O_p(a_n)$  denotes stochastic boundedness with rate  $a_n$  in  $\mu$ . i.e.,  $\mathbf{X} = O_p(a_n)$  (as  $N \rightarrow \infty$ )  $\Leftrightarrow$  for all  $\varepsilon > 0$ , there exist  $\delta(\varepsilon) > 0$  and  $N(\varepsilon) > 0$  such that  $\mu(|\mathbf{X}| / a_n \geq \delta(\varepsilon)) < \varepsilon$  for all  $n \geq N(\varepsilon)$ .

**Preliminary Concepts.**  $P$  and  $Q$  are used as the probability measures on  $(\Omega, \mathcal{F})$ , where  $\mathcal{F}$  denotes the  $\sigma$ -algebra on  $\Omega$ .  $P$  is called *absolutely continuous* with respect to  $Q$ ,  $P(A) = 0$  whenever  $Q(A) = 0$  for any  $A \in \mathcal{F}$ , which is represented as  $P \ll Q$ .  $\frac{dP}{dQ}$  denotes the Radon–Nikodým derivative of  $P$  with respect to  $Q$  for  $P$  and  $Q$  with  $P \ll Q$ .  $\mu$  denotes a probability measure on  $\Omega$  with  $P \ll \mu$  and  $Q \ll \mu$ . An example of  $\mu$  is  $(P + Q)/2$ .  $E_P[\cdot]$  denotes the expectation under the distribution  $P$ , i.e.,  $E_P[\phi(\mathbf{x})] = \int_{\Omega_p} \phi(\mathbf{x}) dP(\mathbf{x})$ , where  $\phi(\mathbf{x})$  represents a measurable function over  $\Omega$ .

**Setup.**  $P$  and  $Q$  are probability distributions on  $\Omega \subset \mathbb{R}^d$  with unknown probability densities  $p$  and  $q$ , respectively. We assume  $p(\mathbf{x}) > 0 \Leftrightarrow q(\mathbf{x}) > 0$  almost everywhere  $\mathbf{x} \in \Omega$ .<sup>1</sup>

## 2.2 DRE WITH $f$ -DIVERGENCE VARIATIONAL REPRESENTATION

Herein, we introduce the  $f$ -divergence variational representation along with the corresponding loss functions used for DRE.

**Definition 2.1** ( $f$ -divergence). The  $f$ -divergence  $D_f$  between two probability measures  $P$  and  $Q$  is induced by a convex function  $f$  that satisfies  $f(1) = 0$ , which can be defined as  $D_f(Q||P) = E_P[f(dQ/dP(\mathbf{x}))]$ .

Various divergences are specific instances derived by choosing an appropriate generator function  $f$ . For example, the function  $f(u) = u \cdot \log u$  yields the Kullback–Leibler divergence.

We then derive the variational representations of  $f$ -divergences using the Legendre transform of the convex conjugate of a twice differentiable convex function  $f$ ,  $f^*(\psi) = \sup_{u \in \mathbb{R}} \{\psi \cdot u - f(u)\}$  (Nguyen et al., 2007):

$$D_f(Q||P) = \sup_{\phi \geq 0} \left\{ E_Q[f'(\phi)] - E_P[f^*(f'(\phi))] \right\}, \quad (1)$$

where the supremum is required over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}$  with  $E_Q[|f'(\phi)|] < \infty$  and  $E_P[|f^*(f'(\phi))|] < \infty$ . The maximum value is achieved at  $\phi(\mathbf{x}) = dQ/dP(\mathbf{x})$ . Pairs of the terms  $f'(\phi)$  and  $f^*(f'(\phi))$  in Equation (1) for major  $f$ -divergences, along with their corresponding convex functions  $f$ , are provided in Table 2 in the Appendix.

By substituting  $\phi$  with a neural network model  $\phi_\theta$  and replacing the expectation  $E$  with sample means  $\hat{E}$ , the optimal function for Equation (1) is trained through back-propagation using an  $f$ -divergence loss function.

$$\mathcal{L}_f(\phi_\theta) = - \left\{ \hat{E}_Q[f'(\phi_\theta)] - \hat{E}_P[f^*(f'(\phi_\theta))] \right\}. \quad (2)$$

Formally, we define the  $f$ -divergence loss function within a probabilistic theoretical framework as follows:

**Definition 2.2** ( $f$ -Divergence Loss). Let  $\hat{\mathbf{X}}_{P[R]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^R\}$ ,  $\mathbf{X}_P^i \stackrel{\text{iid}}{\sim} P$  denote  $R$  i.i.d. random variables from  $P$ , and let  $\hat{\mathbf{X}}_{Q[S]} = \{\mathbf{X}_Q^1, \mathbf{X}_Q^2, \dots, \mathbf{X}_Q^S\}$ ,  $\mathbf{X}_Q^i \stackrel{\text{iid}}{\sim} Q$  denote  $S$  i.i.d. random variables from  $Q$ . Thereafter, for a twice differentiable convex function  $f$ ,  $f$ -divergence loss  $\mathcal{L}_f^{(R,S)}(\cdot)$  is defined as follows:

$$\mathcal{L}_f^{(R,S)}(\phi) = \frac{1}{S} \cdot \sum_{i=1}^S -f'(\phi(\mathbf{X}_Q^i)) + \frac{1}{R} \sum_{i=1}^R f^*(f'(\phi(\mathbf{X}_P^i))), \quad (3)$$

where  $\phi$  denotes a measurable function over  $\Omega$  such that  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$ .

<sup>1</sup>In this study,  $q(\mathbf{x})/p(\mathbf{x})$  is written for  $\frac{dQ}{dP}(\mathbf{x})$  using the Radon–Nikodým density representation for readability.

162 **3 MAIN RESULTS**  
 163

164 The key findings of this study are twofold. First, we establish common upper and lower bounds for  
 165 the  $L_p$  error in DRE by employing variational  $f$ -divergence optimization. Second, we empirically  
 166 investigate the relationship between KL-divergence, data dimension, and the estimation accuracy  
 167 of DRE through variational  $f$ -divergence optimization. Specifically, we discover that the  $L_p$  error  
 168 significantly increases with the rise in KL-divergence when  $p > 1$ , and this increase is exacerbated  
 169 by the magnitude of the order  $p$ .  
 170

171 **3.1 THEORETICAL RESULTS.**  
 172

173 In this study, we outline the assumptions necessary for deriving the upper and lower bounds of  
 174 the DRE. The assumptions are straightforward and primarily involve the consideration of Lipschitz  
 175 continuous estimators. Specifically, we assume the  $L$ -Lipschitz continuity of the energy function of  
 176 the distributions,  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$ .  
 177

178 **Assumption 3.1** (Assumption for the Upper Bound). The following assumption is imposed on the  
 probability distributions  $P$  and  $Q$ .  
 179

180 U1.  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$  is  $L$ -Lipschitz continuous with  $L > 0$  on  $\Omega$ .  
 181

182 **Assumption 3.2** (Assumptions for the Lower Bound). The following assumptions are imposed on  
 the probability distributions  $P$  and  $Q$ .  
 183

184 L1.  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$  is  $L$ -bi-Lipschitz continuous with  $L > 1$  on  $\Omega$ .  
 185

186 L2.  $E_P [(dQ/dP)^p] < \infty$  where  $p \leq d$ .  
 187

188 For the probability distributions  $P$  and  $Q$ , Assumption L1 is crucial for deriving the lower bound of  
 189 the  $L_p$  error in DRE. Further details on this assumption can be found in Remark 4.6 in Section 4.2.  
 190

191 Additionally, Assumptions 3.3 and 3.4 are necessary for deriving both the upper and lower bounds  
 192 of the DRE.  
 193

194 **Assumption 3.3** (Assumptions for the Convex Function  $f$ ). The convex function  $f$  is assumed to  
 195 satisfy the following: (F1)  $f$  is three-times differentiable; (F2)  $f''(u) > 0$  for all  $u > 0$ ; and (F3)  
 196  $E_P [f''(dQ/dP)] < \infty$ .  
 197

198 **Assumption 3.4** (Assumption for the Support). The support  $\Omega$  is assumed to satisfy the following:  
 199 (O1)  $\text{diag}(\Omega) < \infty$ .  
 200

201 Under these conditions, we obtain the upper and lower bounds for the  $L_p$  error in DRE through  
 202 variational  $f$ -divergence optimization.  
 203

204 **Theorem 3.5** (Informal. See Theorem 4.5 and 4.8). Assume  $\Omega$  is a compact set in  $\mathbb{R}^d$ , where  
 205  $d \geq 3$ , and  $f$  satisfies Assumption 3.3. Let  $P$  and  $Q$  denote the probability measures on  $\Omega$ , and  
 206 let  $\phi$  represent a  $K$ -Lipschitz function that minimizes the  $f$ -divergence loss functions. Let  $\phi$  be a  
 207  $K$ -Lipschitz function that minimizes the  $f$ -divergence loss functions  $\mathcal{L}_f^{(R,S)}(\cdot)$  defined as Equation  
 208 (3) using early stopping. Additionally, let  $N = \min\{R, S\}$ .  
 209

210 **(Upper Bound)** Assume Assumption 3.1: Thereafter, Equation (4) holds for  $1 \leq p \leq d/2$  such that  
 211

$$\left\| \frac{q(\mathbf{x})}{p(\mathbf{x})} - \phi(\mathbf{x}) \right\|_{L_p(\Omega, P)} \lesssim \frac{\text{diag}(\Omega)}{N^{1/d}} \cdot \left\{ L \cdot E \left[ \left( \frac{dQ}{dP} \right)^{2 \cdot p} \right]^{1/(2 \cdot p)} + K \right\}. \quad (4)$$

212 **(Lower Bound)** Assume Assumption 3.2: Equations (5) and (6) hold for  $1 \leq p \leq d$  such that  
 213

$$E_{\mathbf{X}_P^1 \dots \mathbf{X}_P^N} \left[ \left\| \frac{q(\mathbf{x})}{p(\mathbf{x})} - \phi(\mathbf{x}) \right\|_{L_p(\Omega, P)} \right] \gtrsim \frac{1}{N^{1/d}} \cdot \left\{ \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - K \cdot \text{diag}(\Omega) \right\} \quad (5)$$

$$\gtrsim \frac{1}{N^{1/d}} \cdot \left\{ \frac{1}{L} \cdot e^{\frac{(p-1)}{p} \cdot KL(P||Q) - 1} - K \cdot \text{diag}(\Omega) \right\}, \quad (6)$$

216 where  $\|\cdot\|_{L_p(\Omega, P)}$  denotes the  $L_p$  norm on  $\Omega$  and the Lebesgue integral on  $P$  and  $KL(P||Q)$  denotes  
 217 the KL-divergence between  $P$  and  $Q$ .  
 218

219 These bounds are applicable to all  $K$ -Lipschitz continuous estimators optimized using the  $f$ -  
 220 divergence loss functions with early stopping, as discussed in Section 4.3 and supported by Theorem  
 221 4.8.

222 Theorem 3.5 indicates that the curse of dimensionality occurs when  $p = 1$ . For  $p > 1$ , both the curse  
 223 of dimensionality and the large sample requirement for high KL-divergence data occur concurrently.  
 224 Equation (6) demonstrates that the  $L_p$  error escalates significantly with increasing KL-divergence  
 225 for  $p > 1$ , and this increase accelerates as  $p$  increases. These theoretical findings are corroborated  
 226 by numerical experiments, which are discussed in the subsequent section.  
 227

### 228 3.2 EXPERIMENTAL RESULTS.

229 We empirically verified the relationship among KL-divergence, data dimension, and estimation ac-  
 230 curacy of DRE through variational  $f$ -divergence optimization. The results, which support the impli-  
 231 cations of Theorem 3.5, are detailed in Section D in the Appendix.  
 232

233  **$L_p$  Errors vs. the KL-Divergence in Data** We conducted the experiments on the relationship  
 234 between  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE and the KL-divergence of the data. In the experiments, we  
 235 generated 100 sets of 5-dimensional datasets with the KL-divergence of 1, 2, 4, 6, 8, 10, 12, and  
 236 14. For each dataset, DRE was conducted using  $\alpha$ -divergence and KL-divergence loss functions,  
 237 then  $L_1$ ,  $L_2$ , and  $L_3$  errors were observed. We reported the results as Figure 1. The details on the  
 238 experimental settings and neural network training are provided in Section D in the Appendix.

239 As displayed in Figure 1, the estimation errors for  $p > 0$  increased significantly, which accelerates  
 240 as  $p$  becomes larger. In contrast, when  $p = 0$ , a relatively mild increase was observed. As indicated  
 241 by Theorem 3.5, these results highlight the impact of the KL-divergence in the data on  $L_p$  error with  
 242  $p > 1$  in DRE  $f$ -divergence loss functions.  
 243

244  **$L_p$  Errors vs. the Dimensions of Data** We conducted experiments to investigate the relationship  
 245 between  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE and the dimensionality of the data. In the experiments, we  
 246 generated 100 sets of datasets of 50, 100 and 200 dimensions with the KL-divergence amounts of  
 247 3. For each dataset, DRE was conducted using  $\alpha$ -divergence and KL-divergence loss functions,  
 248 then  $L_1$ ,  $L_2$ , and  $L_3$  errors were observed. We reported the results as Figure 2. The details on the  
 249 experimental settings and neural network training are provided in Section D in the Appendix.

250 As depicted in Figure 2, the estimation errors  $L_1$ ,  $L_2$ , and  $L_3$  for  $p > 0$  increased as the data  
 251 dimensionality increased for both the  $\alpha$ -divergence and KL-divergence loss functions. These results  
 252 indicate that the curse of dimensionality occurs equally across the  $L_p$  errors, as stated by Theorem  
 253 3.5.  
 254

## 255 4 OVERVIEW OF UPPER AND LOWER BOUND DERIVATIONS

256 In this section, we outline the derivation of the upper and lower bounds. We begin by introducing a  
 257 conceptual reformulation of the  $f$ -divergence loss function, which forms the basis of our theoretical  
 258 framework. Next, we derive the upper and lower bounds for DRE in terms of  $L_P$  error, based on this  
 259 reformulation. Finally, we extend these results to the optimization of the  $f$ -divergence loss function,  
 260 incorporating early stopping and monitoring validation losses, which constitutes the core theoretical  
 261 contribution of this study. Detailed statements and proofs for the theorems mentioned in this section  
 262 are provided in Section C of the Appendix.  
 263

### 264 4.1 CONCEPTUAL REFORMULATION OF THE $f$ -DIVERGENCE LOSS FUNCTIONS

265 The optimization of  $f$ -divergence loss functions, denoted as  $\mathcal{L}_f^{(R,S)}(\phi)$  in Equation (3), presents  
 266 both practical and theoretical challenges owing to their tendency to overfit the training data.  
 267

268 To more deeply understand this issue, let us consider a deterministic setting as described in Defini-  
 269 tion 2.2, where  $(\mathbf{x}_P^1, \mathbf{x}_P^2, \dots, \mathbf{x}_P^R) = (1, 2, \dots, R)$  and  $(\mathbf{x}_Q^1, \mathbf{x}_Q^2, \dots, \mathbf{x}_Q^S) = (R+1, R+2, \dots, R+$

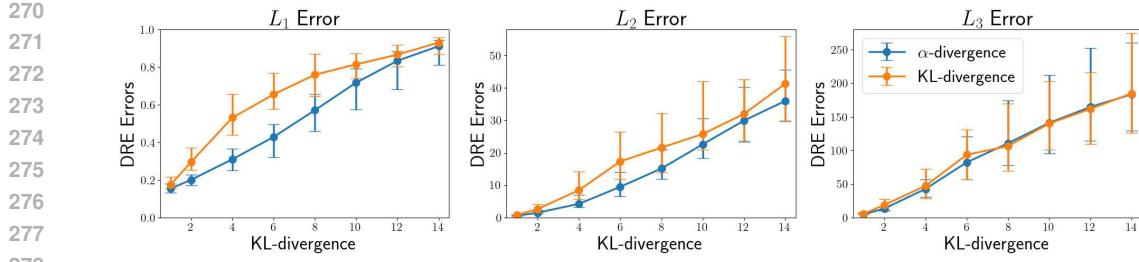


Figure 1: Experimental results of  $L_p$  errors versus the amount of KL-divergence in the data are presented, as detailed in Section 3.2. The  $x$ -axis represents the amount of KL-divergence in synthetic datasets of fixed dimension. The  $y$ -axes of the left, center, and right graphs correspond to the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE, respectively. The plots show the median  $y$ -axis values, while the error bars represent the interquartile range (25th to 75th percentiles). The blue line shows errors using the  $\alpha$ -divergence loss function, and the orange line shows errors using the KL-divergence loss function.

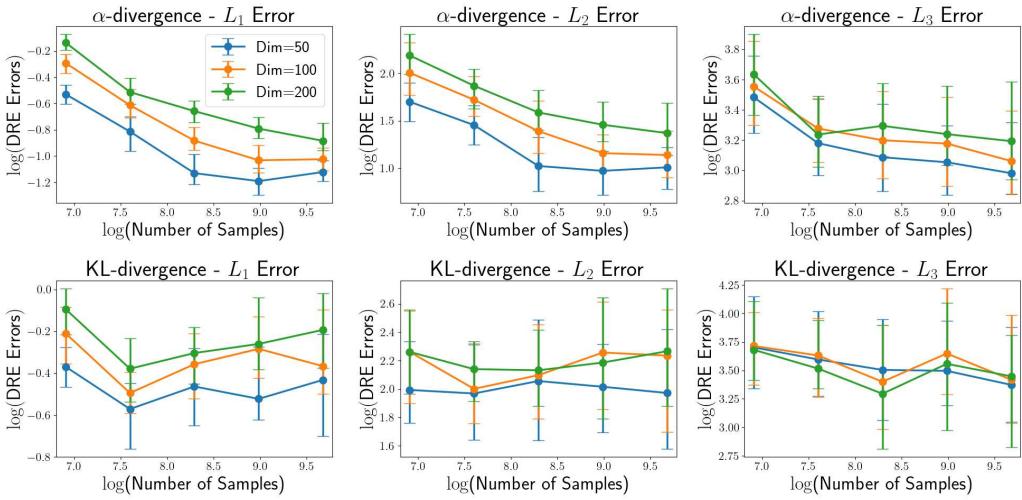


Figure 2: Experimental results on  $L_p$  errors versus the dimensionality of the data are presented, as detailed in Section 3.2. The top row displays results using the  $\alpha$ -divergence loss function, whereas the bottom row presents results using the KL-divergence loss function. The  $x$ -axis represents the logarithm of the number of samples utilized in the optimizations of DRE. The  $y$ -axes of the left, center, and right graphs correspond to the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE, respectively. The plots show the median  $y$ -axis values, while the error bars represent the interquartile range (25th to 75th percentiles). The blue, orange, and green lines show results for data dimensions of 50, 100, and 200, respectively.

$S$ ). Notably,  $\{\mathbf{x}_P^i\}_{i=1}^R \cap \{\mathbf{x}_Q^i\}_{i=1}^S = \emptyset$ . In this setup, we observe that  $\hat{\mathcal{L}}_f^{(R,S)}(\phi) \rightarrow -\infty$  as  $f^*(f'(\phi(\mathbf{x}_P^i))) \rightarrow -\infty$  and  $-f'(\phi(\mathbf{x}_Q^j)) \rightarrow -\infty$  for all  $1 \leq i \leq R$  and  $1 \leq j \leq S$ . In practice, this issue is addressed by implementing early stopping based on monitoring validation losses during optimization. The present theoretical framework accommodates this practical strategy, which facilitates an analysis of both the optimization process and its implications for downstream tasks such as DRE.

To reconcile the practical and theoretical behaviors of  $f$ -divergence loss functions within our framework, we introduce a conceptual reformulation of the loss function.

**Definition 4.1** ( $\mu$ -Representation  $f$ -Divergence Loss). Let  $\mu$  be a probability measure with  $P \ll \mu$  and  $Q \ll \mu$ , and let  $\hat{\mathbf{X}}_{\mu[N]} = \{\mathbf{X}_\mu^1, \dots, \mathbf{X}_\mu^N\}$  denote  $N$  i.i.d. random variables from  $\mu$ . For a twice

324 differentiable convex function  $f$ , let  
 325

$$\tilde{l}_f(u; \mathbf{x}) = -f'(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + f^*(f'(u)) \cdot \frac{dP}{d\mu}(\mathbf{x}), \quad (7)$$

328 where  $f^*$  denotes the Legendre transform of  $f$ :  $f^*(\psi) = \sup_{u \in \mathbb{R}} \{\psi \cdot u - f(u)\}$ . Additionally, let  
 329  $N = \min\{R, S\}$ .

330 The  $\mu$ -representation of the  $f$ -divergence loss  $\mathcal{L}_f^{(R,S)}(\cdot)$  in Equation (3) at the points  $\hat{\mathbf{X}}_{\mu[N]}$  is de-  
 331 fined as  
 332

$$\tilde{\mathcal{L}}_f^{(N)}(\phi) = \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(\phi; \mathbf{X}_\mu^i), \quad (8)$$

333 where  $\phi$  is a measurable function over  $\Omega$  such that  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$ .  
 334

335 This representation introduces an error of  $1/\sqrt{N}$  between the practical  $f$ -divergence loss function  
 336  $\mathcal{L}_f^{(R,S)}(\phi)$  and the  $\mu$ -representation  $f$ -divergence loss  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$ . However, this error is negligible  
 337 when  $d \geq 3$ , which will be discussed in Section 4.3.

338 The optimization properties of this conceptual loss function are encapsulated in Proposition 4.2.  
 339

340 **Proposition 4.2.** *Assume that  $f$  satisfies Assumption 3.3. Let  $\phi_* = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$ .  
 341 Then,  $\phi_*(\mathbf{X}_\mu^i) = \frac{dQ}{dP}(\mathbf{X}_\mu^i)$ , for  $i = 1, 2, \dots, N$ .*

342 This reformulation ensures that the conceptual loss function does not diverge. Furthermore, all  
 343 optimal points in the conceptual loss function are aligned with the ideal density ratios.  
 344

## 345 4.2 DERIVATION OF UPPER AND LOWER BOUNDS FOR OPTIMAL FUNCTIONS OF THE 346 $\mu$ -REPRESENTATION $f$ -DIVERGENCE LOSS FUNCTIONS 347

348 In this section, we derive upper and lower bounds for the  $L_p$  error in DRE for the optimal function  
 349 of  $\mathcal{L}_f^{(N)}(\cdot)$  defined in the previous section, based on the expected distance between the nearest  
 350 neighbors of each  $\mathbf{X}_\mu^i$ ,  $1 \leq N$ .  
 351

352 Hereafter,  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  denotes the nearest neighbor of  $\mathbf{x}$  in  $\hat{\mathbf{X}}_{\mu[N]} = \{\mathbf{X}_\mu^1, \dots, \mathbf{X}_\mu^N\}$ . Specifically,  
 353 define  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  as  $\mathbf{X}_\mu^i$  in  $\hat{\mathbf{X}}_{\mu[N]}$  such that  $\|\mathbf{X}_\mu^l - \mathbf{x}\|_\infty > \|\mathbf{X}_\mu^i - \mathbf{x}\|_\infty$ , for all  $l < i$ , and  $\|\mathbf{X}_\mu^u - \mathbf{x}\|_\infty \geq \|\mathbf{X}_\mu^i - \mathbf{x}\|_\infty$  for all  $u > i$ . As in the previous section, let  $\phi_* = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$ .  
 354

355 As presented in Proposition 4.2, the optimal points of the  $\mu$ -representation  $f$ -divergence loss functions  
 356  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  coincide with the ideal density ratios. This fact provides the following equation,  
 357 serving as the key bridge between the density ratio and its estimation.  
 358

$$\phi_*(\mathbf{X}_\mu^i) = \frac{dQ}{dP}(\mathbf{X}_\mu^i) = \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \quad (9)$$

359 Based on this equation, we can obtain  
 360

$$\left| \phi_*(\mathbf{x})(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p = \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p. \quad (10)$$

361 Using the triangle inequality in the  $L_p$  norm for the density ratios at  $\mathbf{x}$  and its nearest neighbor, we  
 362 obtain  
 363

$$\begin{aligned} & \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right\}^{1/p} - \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p \right\}^{1/p} \\ & \leq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi_*(\mathbf{x}) \right|^p \right\}^{1/p} \\ & \leq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right\}^{1/p} + \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p \right\}^{1/p}. \end{aligned} \quad (11)$$

Assuming the  $L$ -bi-Lipschitz continuity of the energy function of the density ratio,  $T^*(\mathbf{x}) = -\log q(\mathbf{x})/p(\mathbf{x})$ , we yield

$$\begin{aligned} & \frac{1}{L^p} \left( \frac{dQ}{dP}(\mathbf{x}) \right)^p \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p + O_p \left( \frac{1}{N^{1/(2d)}} \right) \\ & \leq \left| \frac{dQ}{dP}(\mathbf{x}) - \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \\ & \leq L^p \cdot \left( \frac{dQ}{dP}(\mathbf{x}) \right)^p \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p + O_p \left( \frac{1}{N^{1/(2d)}} \right). \end{aligned} \quad (12)$$

Additionally, from the  $K$ -Lipschitz continuity of  $\phi_*(\cdot)$  and Equation (9),

$$\left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p = \left| \phi_*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi_*(\mathbf{x}) \right|^p \leq K^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p. \quad (13)$$

Equations (12) and (13) provide the upper and lower bounds of the difference in density ratios between  $\mathbf{x}$  and its nearest neighbor  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  using their distance.

To evaluate the expectation of the distance between  $\mathbf{x}$  and its nearest neighbor  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$ , we present the following theorems: Theorem 4.3 provides an upper bound for the expectation on the right side of Equation (12); Theorem 4.4 establishes a lower bound for the expectation on the left-hand side.

**Theorem 4.3.** Assume that  $\Omega$  is a compact set. Then, for  $1 \leq p \leq d/2$ ,

$$\begin{aligned} & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\ & \leq \text{diag}(\Omega) \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)}. \end{aligned} \quad (14)$$

**Theorem 4.4.** Let  $P$  and  $Q$  be probability measures on a compact set  $\Omega$  in  $\mathbb{R}^d$  with  $d \geq 1$ . Assume that  $P \ll \lambda$  and  $Q \ll \lambda$ , where  $\lambda$  denotes the Lebesgue measure on  $\mathbb{R}^d$ . Let  $p$  be a positive constant such that  $p \geq 1$ . Assume  $E[(dQ/dP)^p] < \infty$ . Then,

$$\begin{aligned} & \lim_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{X}_{P[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right] \right\}^{1/p} \\ & \geq e^{-1} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p}, \end{aligned} \quad (15)$$

where  $E_{\hat{\mathbf{X}}_{P[N]}}[\cdot]$  denotes the expectation over each variable in  $\hat{\mathbf{X}}_{P[N]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^N\}$ .

Notably, using Jensen's inequality on the right-hand side of Equation (15) in Theorem 4.4, the KL-divergence between  $P$  and  $Q$  appears in the lower bound such that

$$\begin{aligned} e^{-1} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} &= e^{-1} \cdot \left\{ E_Q \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{p-1} \right] \right\}^{1/p} \\ &= E_Q \left[ e^{\frac{p-1}{p} \cdot \log \frac{dQ}{dP}(\mathbf{x}) - 1} \right] \\ &\geq e^{E_Q \left[ \frac{p-1}{p} \cdot \log \frac{dQ}{dP}(\mathbf{x}) - 1 \right]} = e^{\frac{p-1}{p} \cdot KL(Q||P) - 1}. \end{aligned} \quad (16)$$

We derive the upper and lower bounds for the  $L_p$  error in DRE for the optimally estimated functions  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ , as stated in Theorem 4.5.

**Theorem 4.5.** Assume  $\Omega$  is a compact set in  $\mathbb{R}^d$  with  $d \geq 3$ , and that  $f$  satisfies Assumption 3.3. Let  $P$  and  $Q$  be probability measures on  $\Omega$ , assuming that  $P \ll \lambda$  and  $Q \ll \lambda$ , where  $\lambda$  denotes the Lebesgue measure on  $\mathbb{R}^d$ . Let  $T^*(\mathbf{x})$  be the energy function of  $dQ/dP(\mathbf{x})$  defined as

432  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$ . Let  $\tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)}$  denote the set of all  $K$ -Lipschitz continuous functions on  $\Omega$   
 433 that minimize  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ . Specifically, define  
 434

$$435 \quad \tilde{\mathcal{F}}^{(N)} = \left\{ \phi_* : \Omega \rightarrow \mathbb{R}_{>0} \mid \tilde{\mathcal{L}}_f^{(N)}(\phi_*) = \min_{\phi} \tilde{\mathcal{L}}_f^{(N)}(\phi) \right\}, \quad (17)$$

437 and

$$438 \quad \mathcal{F}_{K\text{-}Lip} = \left\{ \phi : \Omega \rightarrow \mathbb{R}_{>0} \mid |\phi(\mathbf{y}) - \phi(\mathbf{x})| \leq K \cdot \|\mathbf{y} - \mathbf{x}\|_{\infty} \text{ for all } \mathbf{y}, \mathbf{x} \in \Omega \right\}. \quad (18)$$

440 Subsequently, let  $\tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)} = \tilde{\mathcal{F}}^{(N)} \cap \mathcal{F}_{K\text{-}Lip}$ .

441 **(Upper Bound)** Assume that  $T^*(\mathbf{x})$  satisfies Assumption 3.1. Thereafter, Equation (19) holds for  
 442  $1 \leq p \leq d/2$ , such that for any  $\phi \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)}$ , such that

$$\begin{aligned} 444 \quad & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\ 445 \quad & \leq L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} + K \cdot \text{diag}(\Omega). \end{aligned} \quad (19)$$

446 **(Lower Bound)** Assume that  $T^*(\mathbf{x})$  satisfies Assumption 3.2. Then, Equations (20) and (21) hold  
 447 for any  $\phi \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)}$ , such that

$$\begin{aligned} 448 \quad & \underline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\ 449 \quad & \geq \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - K \cdot \text{diag}(\Omega) \end{aligned} \quad (20)$$

$$450 \quad \geq \frac{1}{L} \cdot e^{\frac{p-1}{p} \cdot KL(Q||P)-1} - K \cdot \text{diag}(\Omega) \quad (21)$$

451 **Remark 4.6.** Equation (12) when  $L = 1$  suggests that  $|dQ/dP(\mathbf{y}) - dQ/dP(\mathbf{x})| = \|\mathbf{y} - \mathbf{x}\|_{\infty}$ ,  
 452 for all  $\mathbf{x}$  and  $\mathbf{y}$  in  $\Omega$ . This typical case is when  $dQ/dP(x_1, x_2, \dots, x_d) \equiv dQ/dP(x_1, x_2, \dots, x_{d'})$   
 453 with  $d' < d$ . Therefore, this case typically occurs when  $dQ/dP(\mathbf{x})$  is a replication of its lower-  
 454 dimensional distribution. In this case, the upper and lower bounds for the  $L_p$  error in DRE are  
 455 considered to follow the lower dimension.

### 465 4.3 DERIVATION OF UPPER AND LOWER BOUNDS FOR OPTIMAL FUNCTIONS OF THE 466 $f$ -DIVERGENCE LOSS FUNCTIONS 467

468 To establish upper and lower bounds for practical DRE using  $f$ -divergence loss function optimization,  
 469 we initially statistically evaluate the discrepancy between the outputs from the practically op-  
 470 timized functions  $\mathcal{L}_f^{(R,S)}(\cdot)$ , employing early stopping based on validation losses, and the theore-  
 471 tically optimized functions  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ . Next, we demonstrate that this discrepancy is negligible when  
 472  $d \geq 3$ . Finally, the upper and lower bounds for DRE are expressed in terms of  $L_p$  error for the  
 473  $f$ -divergence loss function optimization using early stopping, which constitutes the final theoretical  
 474 result of this study.

475 First, according to the central limit theorem, an error of order  $1/\sqrt{N}$  in probability occurs when  
 476 measuring validation losses.

$$477 \quad \mathcal{L}_f^{(R,S)}(\phi) - E_{\mu} [\mathcal{L}_f^{(R,S)}(\phi)] = O_p \left( \frac{1}{\sqrt{N}} \right). \quad (22)$$

478 Equation (22) implies that there is an error margin of  $O_p \left( \frac{1}{\sqrt{N}} \right)$  when monitoring the validation  
 479 losses for early stopping in the optimization of  $\mathcal{L}_f^{(R,S)}(\phi)$ .

480 Subsequently, we utilize the following theorem to demonstrate that the optimization of Equation  
 481 (22), employing early stopping based on validation losses, is governed by the optimization of the  
 482  $\mu$ -representation  $f$ -divergence loss functions  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ .

486     **Theorem 4.7.** Assume the same assumptions as in Proposition 4.2. Let  $\phi_* =$   
 487      $\arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$ . Therefore, for any measurable function  $\phi: \Omega \rightarrow \mathbb{R}_{>0}$ ,  
 488

$$\begin{aligned} 489 \quad & \phi(\mathbf{X}_\mu^i) - \phi_*(\mathbf{X}_\mu^i) = O_p\left(\frac{1}{\sqrt{N}}\right), \quad \text{for } 1 \leq i \leq N, \\ 490 \quad & \iff \mathcal{L}_f^{(R,S)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu [\mathcal{L}_f^{(R,S)}(\phi)] = O_p\left(\frac{1}{\sqrt{N}}\right), \end{aligned} \quad (23)$$

494 where  $\{\mathbf{X}_\mu^1, \mathbf{X}_\mu^2, \dots, \mathbf{X}_\mu^N\}$  is defined in Definition 4.1.  
 495

496 In Equation (23), the first term on the right-hand side denotes the empirical risk of  $\mathcal{L}_f^{(R,S)}(\phi)$   
 497 using validation data, whereas the second term represents the minimum value of its true error.  
 498 This equation illustrates that when  $\mathcal{L}_f^{(R,S)}(\phi)$  is within the actual early stopping margin, specifi-  
 499 cally  $O_p\left(\frac{1}{\sqrt{N}}\right)$ , the function  $\phi$  deviates from the optimal function of  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  by no more than  
 500  $O_p\left(\frac{1}{\sqrt{N}}\right)$ .  
 501

503 Based on Equation (23), we define the optimal function of  $\mathcal{L}_f^{(R,S)}(\phi)$  for use with early stopping  
 504 while monitoring validation losses as follows:  
 505

$$\begin{aligned} 506 \quad & \phi_{\text{val}} \text{ is optimal in the optimization of } \mathcal{L}_f^{(R,S)}(\phi) \text{ using early stopping} \\ 507 \quad & \triangleq \phi_* + O_p\left(\frac{1}{\sqrt{N}}\right), \quad \text{where } \phi_* = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu [\tilde{\mathcal{L}}_f^{(N)}(\phi)]. \end{aligned} \quad (24)$$

510 The difference  $O_p\left(\frac{1}{\sqrt{N}}\right)$ , appearing in Equation (24), is negligible for DRE when  $d \geq 3$ . Indeed,  
 511 using the triangle inequality in the  $L_p$  norm for  $\phi_* = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$  and Equation (20),  
 512 we observe  
 513

$$\begin{aligned} 514 \quad & \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi_{\text{val}}(\mathbf{x}) \right|^p \right\}^{1/p} \geq \underbrace{\left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi^*(\mathbf{x}) \right|^p \right\}^{1/p}}_{=O\left(\frac{1}{N^{1/d}}\right)} - \underbrace{\left\{ E_P \left| \phi_{\text{val}}(\mathbf{x}) - \phi^*(\mathbf{x}) \right|^p \right\}^{1/p}}_{=O\left(\frac{1}{\sqrt{N}}\right)} \ll \frac{1}{N^{1/d}}. \end{aligned} \quad (25)$$

520 Therefore, we finally obtain the following Theorem 4.8.

521 **Theorem 4.8.** Assume the same assumptions and notations as in Theorem 4.5. Additionally, define  
 522

$$523 \quad \mathcal{F}_{K\text{-}Lip}^{(N)} = \left\{ \phi \in \mathcal{F}_{K\text{-}Lip} \mid \exists \phi_* \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)} \text{ such that } \phi = \phi_* + O_p\left(\frac{1}{\sqrt{N}}\right) \right\}. \quad (26)$$

525 That is,  $\mathcal{F}_{K\text{-}Lip}^{(N)}$  denotes the set of all functions that differ by at most  $O_p\left(\frac{1}{\sqrt{N}}\right)$  from some functions  
 526 that minimize  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ . Therefore, the same results as in Theorem 4.5 hold for all  $\phi \in \mathcal{F}_{K\text{-}Lip}^{(N)}$ .  
 527

## 5 CONCLUSIONS

531 We have established upper and lower bounds on the  $L_p$  errors in DRE through the optimization of  
 532  $f$ -divergence loss functions. These bounds are applicable to any member of a group of Lipschitz  
 533 continuous estimators, regardless of the specific  $f$ -divergence loss function used. These bounds pro-  
 534 vide new insights into how the dimensionality of data and the KL divergence between distributions  
 535 affect the accuracy of DRE. Furthermore, the numerical experiments corroborate these theoretical  
 536 findings, demonstrating that the relationship between  $L_p$  errors, KL divergence, and data dimen-  
 537 sionality aligns with the theoretical implications derived from the bounds. This research faces limita-  
 538 tions, particularly in high-dimensional settings where the curse of dimensionality and large sample  
 539 requirements pose challenges. Future studies could refine the theoretical framework to explore loss  
 functions that improve DRE in complex, high-dimensional tasks.

540 REFERENCES  
541

- 542 Anonymous.  $\alpha$ -divergence loss function for neural density ratio estimation (under review), included  
543 in the supplemental materials for this submission. 2024.
- 544 Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial  
545 networks. *arXiv preprint arXiv:1701.04862*, 2017.
- 546 Francis Bach. Self-concordant analysis for logistic regression. 2010.
- 548 Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeshwar, Sherjil Ozair, Yoshua Bengio, Aaron  
549 Courville, and Devon Hjelm. Mutual information neural estimation. In *International conference  
550 on machine learning*, pp. 531–540. PMLR, 2018.
- 551 Denis Belomestny, Eric Moulines, Alexey Naumov, Nikita Puchkin, and Sergey Samsonov. Rates  
552 of convergence for density estimation with gans. *arXiv preprint arXiv:2102.00199*, 2021.
- 554 Gérard Biau and Luc Devroye. *Lectures on the nearest neighbor method*, volume 246. Springer,  
555 2015.
- 556 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,  
557 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information  
558 processing systems*, 27, 2014.
- 559 Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle  
560 for unnormalized statistical models. In *Proceedings of the thirteenth international conference on  
561 artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings,  
562 2010.
- 564 R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam  
565 Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation  
566 and maximization. *arXiv preprint arXiv:1808.06670*, 2018.
- 567 Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola. Correct-  
568 ing sample selection bias by unlabeled data. *Advances in neural information processing systems*,  
569 19, 2006.
- 571 Takafumi Kanamori, Taiji Suzuki, and Masashi Sugiyama. Statistical analysis of kernel-based least-  
572 squares density-ratio estimation. *Machine Learning*, 86:335–367, 2012.
- 573 Masahiro Kato and Takeshi Teshima. Non-negative bregman divergence minimization for deep  
574 direct density ratio estimation. In *International Conference on Machine Learning*, pp. 5320–  
575 5333. PMLR, 2021.
- 577 Masanari Kimura and Howard Bondell. Density ratio estimation via sampling along generalized  
578 geodesics on statistical manifolds. *arXiv preprint arXiv:2406.18806*, 2024.
- 579 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
580 2014.
- 582 Samory Kpotufe. Lipschitz density-ratios, structured data, and data-driven tuning. In *Artificial  
583 Intelligence and Statistics*, pp. 1320–1328. PMLR, 2017.
- 584 Tengyuan Liang. How well can generative adversarial networks learn densities: A nonparametric  
585 view. *arXiv preprint arXiv:1712.08244*, 2017.
- 587 Zhexiao Lin, Peng Ding, and Fang Han. Estimation based on nearest neighbor matching: from  
588 density ratio to average treatment effect. *Econometrica*, 91(6):2187–2217, 2023.
- 589 David McAllester and Karl Stratos. Formal limitations on the measurement of mutual information.  
590 In *International Conference on Artificial Intelligence and Statistics*, pp. 875–884. PMLR, 2020.
- 591 XuanLong Nguyen, Martin J Wainwright, and Michael Jordan. Estimating divergence functionals  
592 and the likelihood ratio by penalized convex risk minimization. *Advances in neural information  
593 processing systems*, 20, 2007.

- 594 XuanLong Nguyen, Martin J Wainwright, and Michael I Jordan. Estimating divergence functionals  
 595 and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*,  
 596 56(11):5847–5861, 2010.
- 597
- 598 Jonathan Niles-Weed and Quentin Berthet. Minimax estimation of smooth densities in wasserstein  
 599 distance. *The Annals of Statistics*, 50(3):1519–1540, 2022.
- 600 Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. f-gan: Training generative neural samplers  
 601 using variational divergence minimization. *Advances in neural information processing systems*,  
 602 29, 2016.
- 603
- 604 George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lak-  
 605 shminarayanan. Normalizing flows for probabilistic modeling and inference. *Journal of Machine  
 606 Learning Research*, 22(57):1–64, 2021.
- 607 Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,  
 608 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in  
 609 pytorch. 2017.
- 610 Ben Poole, Sherjil Ozair, Aaron Van Den Oord, Alex Alemi, and George Tucker. On variational  
 611 bounds of mutual information. In *International Conference on Machine Learning*, pp. 5171–  
 612 5180. PMLR, 2019.
- 613
- 614 Nikita Puchkin, Sergey Samsonov, Denis Belomestny, Eric Moulines, and Alexey Naumov. Rates  
 615 of convergence for density estimation with generative adversarial networks. *Journal of Machine  
 616 Learning Research*, 25(29):1–47, 2024.
- 617 Igal Sason and Sergio Verdú.  $f$ -divergence inequalities. *IEEE Transactions on Information Theory*,  
 618 62(11):5973–6006, 2016.
- 619
- 620 Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-  
 621 likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 622 Shashank Singh and Barnabás Póczos. Minimax distribution estimation in wasserstein distance.  
 623 *arXiv preprint arXiv:1802.08855*, 2018.
- 624
- 625 Jiaming Song and Stefano Ermon. Understanding the limitations of variational mutual information  
 626 estimators. *arXiv preprint arXiv:1910.06222*, 2019.
- 627 Charles J Stone. Optimal rates of convergence for nonparametric estimators. *The annals of Statistics*,  
 628 pp. 1348–1360, 1980.
- 629
- 630 Masashi Sugiyama, Taiji Suzuki, and Takafumi Kanamori. Density-ratio matching under the breg-  
 631 man divergence: a unified framework of density-ratio estimation. *Annals of the Institute of Sta-  
 632 tistical Mathematics*, 64:1009–1044, 2012.
- 633 Masatoshi Uehara, Issei Sato, Masahiro Suzuki, Kotaro Nakayama, and Yutaka Matsuo. Generative  
 634 adversarial nets from a density ratio estimation perspective. *arXiv preprint arXiv:1610.02920*,  
 635 2016.
- 636
- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647

---

## 648 A ORGANIZATION OF THE SUPPLEMENTARY DOCUMENT 649

650 The organization of this supplementary document is as follows: Section B provides a list of notations  
651 used in this study. Section C presents the proofs referenced in Sections 3 and 4. Section D provides  
652 details of the experiments conducted. Section E explores further discussions related to this study.  
653

654 Additionally, the code used in the numerical experiments is included as supplementary material.  
655

## 656 B NOTATIONS 657

658 We list all notations used in the Appendix used in this study in Table 1.  
659

## 660 C PROOFS 661

662 In this section, we present the theorems and proofs referenced in this study. We begin by summarizing  
663 all definitions and assumptions stated in previous sections. Then, we provide the theorems and  
664 proofs used throughout this study.  
665

### 666 C.1 DEFINITIONS AND ASSUMPTIONS IN SECTIONS 2, 3, AND 4 667

#### 668 C.1.1 DEFINITIONS 669

670 **Definition C.1** ( $f$ -Divergence (Definition 2.1 restated)). The  $f$ -divergence  $D_f$  between two probability  
671 measures  $P$  and  $Q$ , which is induced by a convex function  $f$  satisfying  $f(1) = 0$ , is defined as  
672  $D_f(Q||P) = E_P[f(q(\mathbf{x})/p(\mathbf{x}))]$ .  
673

674 **Definition C.2** ( $f$ -Divergence Loss (Definition 2.2 restated)). Let  $\hat{\mathbf{X}}_{P[R]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^R\}$ ,  
675  $\mathbf{X}_P^i \stackrel{\text{iid}}{\sim} P$  denote  $R$  i.i.d. random variables from  $P$ , and let  $\hat{\mathbf{X}}_{Q[S]} = \{\mathbf{X}_Q^1, \mathbf{X}_Q^2, \dots, \mathbf{X}_Q^S\}$ ,  $\mathbf{X}_Q^i \stackrel{\text{iid}}{\sim} Q$   
676 denote  $S$  i.i.d. random variables from  $Q$ . Then, for a twice differentiable convex function  $f$ ,  $f$ -divergence loss  
677  $\mathcal{L}_f^{(R,S)}(\cdot)$  is defined as follows:  
678

$$679 \mathcal{L}_f^{(R,S)}(\phi) = \frac{1}{S} \cdot \sum_{i=1}^S -f'(\phi(\mathbf{X}_Q^i)) + \frac{1}{R} \sum_{i=1}^R f^*(f'(\phi(\mathbf{X}_P^i))), \quad (27)$$

680 where  $\phi$  is a measurable function over  $\Omega$  such that  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$ .  
681

682 **Definition C.3** ( $\mu$ -Representation  $f$ -Divergence Loss (Definition 4.1 restated)). Let  $f$  be a twice  
683 differentiable convex function  $f$ . Then,  $\mu$ -representation function of  $f$  for  $u > 0$  at a point  $\mathbf{x} \in \Omega$ ,  
684 which is written for  $\tilde{l}_f(u)$  in an abbreviated form, is defined as  
685

$$686 \tilde{l}_f(u; \mathbf{x}) = -f'(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + f^*(f'(u)) \cdot \frac{dP}{d\mu}(\mathbf{x}), \quad (28)$$

687 where  $f^*$  denotes the Legendre transform of  $f$ :  $f^*(\psi) = \sup_{u \in \mathbb{R}} \{\psi \cdot u - f(u)\}$ . Let  $N =$   
688  $\min\{R, S\}$ , and let  $\hat{\mathbf{X}}_{\mu[N]} = \{\mathbf{X}_{\mu}^1, \dots, \mathbf{X}_{\mu}^N\}$  denote  $N$  i.i.d. random variables from  $\mu$ . Then,  
689  $\mu$ -representation of the  $f$ -divergence loss  $\mathcal{L}_f^{(R,S)}(\cdot)$  in Equation (??) at the points  $\hat{\mathbf{X}}_{\mu[N]}$  is defined  
690 as  
691

$$692 \mathcal{L}_f^{(N)}(\phi) = \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(u; \mathbf{X}_{\mu}^i) \quad (29)$$

693 where  $\phi$  is a measurable function over  $\Omega$  such that  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$ .  
694

#### 695 C.1.2 ASSUMPTIONS 696

697 **Assumption C.4** (Assumption for the Upper Bound (Assumption 3.1 restated)). The following  
698 assumption is imposed on the probability distributions  $P$  and  $Q$ .  
699

702  
703  
704  
705

Table 1: Notations and definitions used in the proofs

Notations	Definitions, Meanings
(Capital, small, and bold letters)	Random variables are denoted by capital letters; for example, $A$ . Small letters are used for values of the random variables of the corresponding capital letters. Bold letters $\mathbf{A}$ and $\mathbf{a}$ represent sets of random variables and their values.
$\mathbb{R}, \mathbb{R}^d$	The set of all real numbers and the $d$ -dimensional vector space over the real numbers, respectively.
$\mathbb{R}_{>0}$	The set of all positive real numbers: $\mathbb{R}_{>0} = \{x \in \mathbb{R} \mid x > 0\}$ .
$\Omega$	A subset of $\mathbb{R}^d$ : $\Omega \subset \mathbb{R}^d$ .
$f(x) = O(g(x))$ , as $x \rightarrow a$	Asymptotic boundedness with rate $g(x)$ as $x \rightarrow a$ : $f(x) = O(g(x)) \Leftrightarrow \limsup_{x \rightarrow a}  f(x)/g(x)  \leq C$ , where $C > 0$ .
$f(x) = o(g(x))$ , as $x \rightarrow a$	Asymptotic domination with rate $g(x)$ as $x \rightarrow a$ : $f(x) = o(g(x)) \Leftrightarrow \lim_{x \rightarrow a} f(x)/g(x) = 0$ .
$\mathbf{X} = O_p(a_N)$ , as $N \rightarrow \infty$	Stochastic boundedness with rate $a_N$ in $\mu$ : $\mathbf{X} = O_p(a_N) \Leftrightarrow$ for all $\varepsilon > 0$ , there exist $\delta(\varepsilon) > 0$ and $N(\varepsilon) > 0$ such that $\mu( \mathbf{X} /a_N \geq \delta(\varepsilon)) < \varepsilon$ for all $N \geq N(\varepsilon)$ .
$\mathbf{X} = o_p(a_N)$ , as $N \rightarrow \infty$	Convergence in probability with rate $a_N$ in $\mu$ : $\mathbf{X} = o_p(a_N) \Leftrightarrow$ for all $\varepsilon > 0$ , for all $\delta > 0$ , there exists $N(\varepsilon, \delta) > 0$ such that $\mu( \mathbf{X} /a_N \geq \delta) < \varepsilon$ for all $N \geq N(\varepsilon)$ .
$P \ll Q$	$P$ is absolutely continuous with respect to $Q$ .
$P, Q$	A pair of probability measures with $P \ll Q$ and $Q \ll P$ .
$\mu$	A probability measure with $P \ll \mu$ and $Q \ll \mu$ .
$\frac{dP}{dQ}$	The Radon–Nikodým derivative of $P$ with respect to $Q$ .
$\hat{\mathbf{X}}_{P[R]}$	$R$ i.i.d. random variables from $P$ : $\hat{\mathbf{X}}_{P[R]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^R\}$ , where $\mathbf{X}_P^i \stackrel{\text{iid}}{\sim} P$ .
$\hat{\mathbf{X}}_{Q[S]}$	$S$ i.i.d. random variables from $Q$ : $\hat{\mathbf{X}}_{Q[S]} = \{\mathbf{X}_Q^1, \mathbf{X}_Q^2, \dots, \mathbf{X}_Q^S\}$ , where $\mathbf{X}_Q^i \stackrel{\text{iid}}{\sim} Q$ .
$N$	$N = \min\{R, S\}$ .
$\hat{\mathbf{X}}_{\mu[N]}$	$N$ i.i.d. random variables from $\mu$ : $\hat{\mathbf{X}}_{\mu[N]} = \{\mathbf{X}_\mu^1, \mathbf{X}_\mu^2, \dots, \mathbf{X}_\mu^N\}$ , where $\mathbf{X}_\mu^i \stackrel{\text{iid}}{\sim} \mu$ .
$\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$	The nearest neighbor variable of $\mathbf{x}$ in $\hat{\mathbf{X}}_{\mu[N]}$ : $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$ is the $\mathbf{X}_\mu^i$ such that $\ \mathbf{X}_\mu^i - \mathbf{x}\  < \ \mathbf{X}_\mu^j - \mathbf{x}\ $ for all $j \neq i$ .
$D_f(Q  P)$	$f$ -divergence: $D_f(Q  P) = E_P[f(q(\mathbf{x})/p(\mathbf{x}))]$ . See Definition C.1.
$\mathcal{L}_f^{(R,S)}(\cdot)$	$f$ -divergence loss function. See Definition C.2.
$\tilde{l}_f(u; \mathbf{x})$	$\mu$ -representation of the $f$ -divergence loss function at $\mathbf{x}$ : $\tilde{l}_f(u; \mathbf{x}) = -f'(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + f^*(f'(u)) \cdot \frac{dP}{d\mu}(\mathbf{x})$ .
$\tilde{\mathcal{L}}_f^{(N)}(\cdot)$	$\mu$ -representation of the $f$ -divergence loss function $\mathcal{L}_f^{(R,S)}(\cdot)$ . See Definition 4.1.
$\bar{\mathcal{L}}_f(\phi)$	The expectation of the $\mu$ -representation of the $f$ -divergence loss on $\mu$ . See Lemma C.11.
$\ \cdot\ $	The Euclidean norm.
$\ \cdot\ _\infty$	The maximum norm in $\mathbb{R}^d$ : $\ \mathbf{y} - \mathbf{x}\ _\infty = \max_{1 \leq i \leq d}  y_i - x_i $ .
$\Delta(\mathbf{a}, r)$	The $d$ -dimensional interval centered at $\mathbf{a}$ with each side of length $r$ : $\Delta(\mathbf{a}, r) = \{\mathbf{x} \in \mathbb{R}^d \mid \ \mathbf{x} - \mathbf{a}\ _\infty < r/2\}$ .
$\text{diag}(\mathcal{B})$	The diameter of $\mathcal{B}$ : $\text{diag}(\mathcal{B}) = \inf_{r \in \mathbb{R}} \{\mathcal{B} \subseteq \Delta(\mathbf{a}, r) \mid \exists \mathbf{a} \in \mathcal{B}\}$ .

754  
755

756    U1.  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$  is  $L$ -Lipschitz continuous with  $L > 0$  on  $\Omega$ . i.e.,  $\exists L > 0$  s.t.  
 757     $|T^*(\mathbf{y}) - T^*(\mathbf{x})| \leq L \cdot \|\mathbf{y} - \mathbf{x}\|_\infty$  for any  $\mathbf{y}, \mathbf{x} \in \Omega$ .

759    **Assumption C.5** (Assumptions for the Lower Bound (Assumption 3.2 restated)). The following  
 760    assumptions are imposed on the probability distributions  $P$  and  $Q$ .

761    L1.  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$  is  $L$ -bi-Lipschitz continuous on  $\Omega$ . i.e.,  $T^*(\mathbf{x}) =$   
 762     $-\log q(\mathbf{x})/p(\mathbf{x})$  is  $L$ -Lipschitz continuous with  $L > 0$  on  $\Omega$ .

764    L2.  $E_P[(dQ/dP)^p] < \infty$  where  $p \leq d$ .

765    **Assumption C.6** (Assumptions for the Convex Function  $f$  (Assumption 3.3 restated)). The following  
 766    assumptions are assumed for the convex function  $f$ .

768    F1.  $f$  is three-time differentiable.

769    F2.  $f''(u) > 0$  for all  $u > 0$ .

771    F3.  $E_P[f''(dQ/dP)] < \infty$ .

772    **Assumption C.7** (Assumption for the Support (Assumption 3.4 restated)). The following assumption  
 773    is assumed for  $\Omega$ .

775    O1.  $\text{diag}(\Omega) < \infty$ .

## 777    C.2 THEOREMS AND PROOFS IN SECTIONS 2, 3, AND 4

779    **Lemma C.8.** Let  $f$  be a twice differentiable function. Consider  $\tilde{l}_f(u; \mathbf{x})$  defined as in Equation  
 780    (28). Then, the first derivative of  $\tilde{l}_f(u; \mathbf{x})$  with respect to  $u$  is given by:

$$782 \quad \frac{d}{du} \tilde{l}_f(u; \mathbf{x}) = \left\{ u - \frac{dQ}{dP}(\mathbf{x}) \right\} \cdot f''(u) \cdot \frac{dP}{d\mu}(\mathbf{x}). \quad (30)$$

784    Additionally, if  $\tilde{l}_f(u; \mathbf{x})$  is thrice differentiable, the second derivative with respect to  $u$  is given by:

$$786 \quad \frac{d^2}{du^2} \tilde{l}_f(u; \mathbf{x}) = \left\{ \left( u - \frac{dQ}{dP}(\mathbf{x}) \right) \cdot f'''(u) + f''(u) \right\} \cdot \frac{dP}{d\mu}(\mathbf{x}). \quad (31)$$

789    *Proof of Lemma C.8.* First, note that

$$791 \quad \begin{aligned} \tilde{l}_f(u; \mathbf{x}) &= -f'(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + f^*(f'(u)) \cdot \frac{dP}{d\mu}(\mathbf{x}) \\ 793 &= -f'(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + \{f'(u) \cdot u - f(u)\} \cdot \frac{dP}{d\mu}(\mathbf{x}). \end{aligned} \quad (32)$$

795    Differentiating Equation (32) with respect to  $u$ , we obtain the first and second derivatives of  $\tilde{l}_f(u; \mathbf{x})$   
 796    as follows:

$$798 \quad \begin{aligned} \frac{d}{du} \tilde{l}_f(u; \mathbf{x}) &= -f''(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + u \cdot f''(u) \cdot \frac{dP}{d\mu}(\mathbf{x}) \\ 800 &= \left\{ u - \frac{dQ}{dP}(\mathbf{x}) \right\} \cdot f''(u) \cdot \frac{dP}{d\mu}(\mathbf{x}), \end{aligned} \quad (33)$$

803    and

$$804 \quad \begin{aligned} \frac{d^2}{du^2} \tilde{l}_f(u; \mathbf{x}) &= -f'''(u) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + f''(u) \cdot \frac{dP}{d\mu}(\mathbf{x}) + u \cdot f'''(u) \cdot \frac{dP}{d\mu}(\mathbf{x}) \\ 806 &= \left\{ \left( u - \frac{dQ}{dP}(\mathbf{x}) \right) \cdot f'''(u) + f''(u) \right\} \cdot \frac{dP}{d\mu}(\mathbf{x}). \end{aligned} \quad (34)$$

809    This completes the proof.  $\square$

**Theorem C.9.** Assume that  $f$  satisfies Assumption C.6. Then,  $\tilde{l}_f(u; \mathbf{x})$ , as defined in Equation (28), is minimized only when  $u^*(\mathbf{x}) = \frac{dQ}{dP}(\mathbf{x})$ . In addition, for  $u > 0$ , the following holds:

$$\begin{aligned} \tilde{l}_f(u; \mathbf{x}) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) \\ = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}) \cdot \left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2 + o\left(\left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2\right), \end{aligned} \quad (35)$$

where  $f(a) = o(a)$  (as  $a \rightarrow 0$ ) denotes asymptotic domination such that  $\lim_{a \rightarrow 0} \frac{f(a)}{a} \rightarrow 0$ .

*Proof of Theorem C.9.* Let  $\text{sign}(x)$  denote the sign of the value  $x$ : specifically,  $\text{sign}(x) = 1$  if  $x > 0$ ,  $\text{sign}(x) = -1$  if  $x < 0$ , and  $\text{sign}(x) = 0$  if  $x = 0$ .

From Equation (30) in Lemma C.8, we have

$$\begin{aligned} \text{sign}\left(\frac{d}{du}\tilde{l}_f(u; \mathbf{x})\right) &= \text{sign}\left(\left\{u - \frac{dQ}{dP}(\mathbf{x})\right\} \cdot f''(u) \cdot \frac{dP}{d\mu}(\mathbf{x})\right) \\ &= \text{sign}\left(\left\{u - \frac{dQ}{dP}(\mathbf{x})\right\}\right) \cdot \text{sign}(f''(u)) \cdot \text{sign}\left(\frac{dP}{d\mu}(\mathbf{x})\right) \\ &= \text{sign}\left(u - \frac{dQ}{dP}(\mathbf{x})\right). \end{aligned} \quad (36)$$

Thus,  $\tilde{l}_f(u; \mathbf{x})$  is minimized only when  $u^* = \frac{dQ}{dP}(\mathbf{x})$ .

Next, from Equation (30),

$$\frac{d}{du}\tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) = 0, \quad (37)$$

and from Equation (31),

$$\frac{d^2}{du^2}\tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) = f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}). \quad (38)$$

Thus, using the second-order Taylor expansion of  $\tilde{l}_f(u; \mathbf{x})$  around  $u = \frac{dQ}{dP}(\mathbf{x})$ , we have

$$\begin{aligned} \tilde{l}_f(u; \mathbf{x}) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) \\ = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}) \cdot \left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2 + o\left(\left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2\right). \end{aligned} \quad (39)$$

This completes the proof.  $\square$

**Proposition C.10** (Proposition 4.2 restated). Assume that  $f$  satisfies Assumption C.6. Let  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  denote the  $\mu$ -representation  $f$ -divergence loss as defined in Definition C.3. Then, the minimum value of  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  is achieved if and only if  $\phi$  satisfies

$$\phi(\mathbf{X}_\mu^i) = \frac{dQ}{dP}(\mathbf{X}_\mu^i), \quad \text{for } i = 1, 2, \dots, N. \quad (40)$$

*proof of Proposition C.10.* From Theorem C.9, we observe that, for  $i = 1, 2, \dots, N$ ,

$$\min_{u>0} \tilde{l}_f(u; \mathbf{X}_\mu^i) = \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right), \quad (41)$$

where the minimum value is archived only at  $u = \frac{dQ}{dP}(\mathbf{X}_\mu^i)$ .

864 Thus,  
865

$$\begin{aligned}
 866 \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) &= \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) \\
 867 &= \min_{\substack{\phi(\mathbf{X}_\mu^i) > 0, \\ i=1,2,\dots,N}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) \\
 868 &= \min_{\substack{u_i > 0, \\ i=1,2,\dots,N}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(u_i; \mathbf{X}_\mu^i) \\
 869 &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right). \tag{42}
 \end{aligned}$$

870 Suppose that  $\tilde{\phi}(\mathbf{x})$  is a function on  $\Omega$  that satisfies Equation (40), we have, from Equation (42),  
871

$$\begin{aligned}
 872 \tilde{\mathcal{L}}_f^{(N)}(\tilde{\phi}) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \\
 873 &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \tilde{\phi}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) \\
 874 &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) \\
 875 &= 0. \tag{43}
 \end{aligned}$$

876 Here, we show that the minimum value of  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  is  
877 archived if  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  satisfies Equation (40).

878 Next, we show that the minimum value of  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  is  
879 archived only if  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  satisfies Equation (40).

880 We have, for any function  $\phi : \Omega \rightarrow (0, \infty)$ ,

$$\begin{aligned}
 881 \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \\
 882 &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f \left( \phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \frac{1}{N} \cdot \sum_{i=1}^N \min_{\substack{u_i > 0, \\ i=1,2,\dots,N}} \tilde{l}_f(u_i; \mathbf{X}_\mu^i) \\
 883 &= \frac{1}{N} \cdot \sum_{i=1}^N \left\{ \tilde{l}_f \left( \phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \min_{u>0} \tilde{l}_f(u; \mathbf{X}_\mu^i) \right\}. \tag{44}
 \end{aligned}$$

884 Suppose that  $\phi(\mathbf{X}_\mu^i) \neq \frac{dQ}{dP}(\mathbf{X}_\mu^i)$ . Then, from Equation (41), we have  
885

$$\tilde{l}_f \left( \phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) > \min_{u>0} \tilde{l}_f(u; \mathbf{X}_\mu^i). \tag{45}$$

886 From Equations (44) and (45), we observe that  
887

$$\begin{aligned}
 888 \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \\
 889 &= \frac{1}{N} \cdot \sum_{i=1}^N \left\{ \tilde{l}_f \left( \phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \min_{u>0} \tilde{l}_f(u; \mathbf{X}_\mu^i) \right\} \\
 890 &\geq \frac{1}{N} \cdot \left\{ \tilde{l}_f \left( \phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i \right) - \min_{u>0} \tilde{l}_f(u; \mathbf{X}_\mu^i) \right\} \\
 891 &> 0 \tag{46}
 \end{aligned}$$

918 Thus, we see that the minimum value of  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  is  
 919 archived only if  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  satisfies Equation (40).  
 920

921 This completes the proof.  $\square$

922 **Lemma C.11.** Assume that  $f$  satisfies Assumption C.6. Let  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  denote the  $\mu$ -representation  
 923  $f$ -divergence loss as defined in Definition C.3. Define  
 924

$$\begin{aligned} 925 \quad \tilde{\mathcal{L}}_f(\phi) &= E_\mu \left[ \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] \\ 926 \\ 927 \quad &= \frac{1}{N} \sum_{i=1}^N E_\mu \left[ -f'(\phi(\mathbf{x}_i)) \cdot \frac{dQ}{d\mu}(\mathbf{x}_i) \right] \\ 928 \\ 929 \quad &\quad + \frac{1}{N} \sum_{i=1}^N E_\mu \left[ f^*(f'(\phi(\mathbf{x}_i))) \cdot \frac{dP}{d\mu}(\mathbf{x}_i) \right]. \end{aligned} \quad (47)$$

930 Then,  
 931

$$932 \quad E_\mu \left[ \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] = \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f(\phi) = \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu \left[ \mathcal{L}_f^{(R,S)}(\phi) \right], \quad (48)$$

933 where the infimum are taken over all measurable functions  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$  such that  $E_P[f(\phi(\mathbf{X}))] <$   
 934  $\infty$ . Additionally, the equality in Equation (48) hold when  $\phi(\mathbf{x}) = \frac{dQ}{dP}(\mathbf{x})$ .  
 935

936 proof of Lemma C.11. Let,  $\tilde{l}_f^*(\mathbf{x}) = \min_{u \in \mathbb{R}_{>0}} \tilde{l}_f(u; \mathbf{x})$ . From Theorem C.9, we see  $\tilde{l}_f^*(\mathbf{x}) =$   
 937  $\tilde{l}_f(dQ/dP(\mathbf{x}); \mathbf{x})$ . Then, we have  
 938

$$\begin{aligned} 939 \quad \tilde{l}_f^*(\mathbf{x}) &= \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{x}); \mathbf{x} \right) \\ 940 \\ 941 \quad &= -f' \left( \frac{dQ}{dP}(\mathbf{x}) \right) \cdot \frac{dQ}{d\mu}(\mathbf{x}) + \left\{ f' \left( \frac{dQ}{dP}(\mathbf{x}) \right) \cdot \frac{dQ}{dP}(\mathbf{x}) - f \left( \frac{dQ}{dP}(\mathbf{x}) \right) \right\} \cdot \frac{dP}{d\mu}(\mathbf{x}) \\ 942 \\ 943 \quad &= -f \left( \frac{dQ}{dP}(\mathbf{x}) \right) \frac{dP}{d\mu}(\mathbf{x}). \end{aligned} \quad (49)$$

944 Now, we have  
 945

$$\begin{aligned} 946 \quad \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) &= \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) \\ 947 \\ 948 \quad &= \min_{\substack{\phi(\mathbf{X}_\mu^i) > 0, \\ i=1,2,\dots,N}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) \\ 949 \\ 950 \quad &= \min_{\substack{u_i > 0, \\ i=1,2,\dots,N}} \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f(u_i; \mathbf{X}_\mu^i) \\ 951 \\ 952 \quad &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f^*(\mathbf{X}_\mu^i). \end{aligned} \quad (50)$$

953 Additionally, we have  
 954

$$\begin{aligned} 955 \quad E_\mu \left[ \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] &= E_\mu \left[ \frac{1}{N} \cdot \sum_{i=1}^N -f'(\phi(\mathbf{x}_i)) \cdot \frac{dQ}{d\mu}(\mathbf{x}_i) \right. \\ 956 \\ 957 \quad &\quad \left. + \frac{1}{N} \cdot \sum_{i=1}^N f^*(f'(\phi(\mathbf{x}_i))) \cdot \frac{dP}{d\mu}(\mathbf{x}_i) \right] \end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{N} \cdot \sum_{i=1}^N E_\mu \left[ f'(\phi(\mathbf{x}_i)) \cdot \frac{dQ}{d\mu}(\mathbf{x}_i) \right] \\
&\quad + \frac{1}{N} \cdot \sum_{i=1}^N E_\mu \left[ f^*(f'(\phi(\mathbf{x}_i))) \cdot \frac{dP}{d\mu}(\mathbf{x}_i) \right] \\
&= -\frac{1}{N} \cdot \sum_{i=1}^N E_Q [f'(\phi)] + \frac{1}{N} \cdot \sum_{i=1}^N E_P [f^*(f'(\phi))] \\
&= -E_Q [f'(\phi)] + E_P [f^*(f'(\phi))], \tag{51}
\end{aligned}$$

and

$$\begin{aligned}
E \left[ \mathcal{L}_f^{(R,S)}(\phi) \right] &= E \left[ \frac{1}{R} \cdot \sum_{i=1}^S -f'(\phi(\mathbf{x}_i^q)) \right. \\
&\quad \left. + \frac{1}{S} \cdot \sum_{i=1}^R f^*(f'(\phi(\mathbf{x}_i^p))) \right] \\
&= -\frac{1}{S} \cdot \sum_{i=1}^S E_Q [f'(\phi(\mathbf{x}_i))] \\
&\quad + \frac{1}{R} \cdot \sum_{i=1}^R E_P [f^*(f'(\phi(\mathbf{x}_i)))] \\
&= -\frac{1}{S} \cdot \sum_{i=1}^S E_Q [f'(\phi)] + \frac{1}{R} \cdot \sum_{i=1}^R E_P [f^*(f'(\phi))] \\
&= -E_Q [f'(\phi)] + E_P [f^*(f'(\phi))]. \tag{52}
\end{aligned}$$

Now, note that, from Equation (1) (Nguyen et al. (2007)), we see

$$\min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} -E_Q [f'(\phi)] + E_P [f^*(f'(\phi))] = -D_f(Q||P), \tag{53}$$

where  $D_f(Q||P)$  denotes  $f$ -divergence defined in Definition C.1 and the equality in Equation (53) holds for  $\phi(\mathbf{x}) = dQ/dP(\mathbf{x})$ .

From Equations (51), (52) and (53), we have

$$\min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu \left[ \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] = \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E \left[ \mathcal{L}_f^{(R,S)}(\phi) \right] = -D_f(Q||P), \tag{54}$$

and the equality in Equation (54) holds for  $\phi(\mathbf{x}) = dQ/dP(\mathbf{x})$ .

Substituting Equation (49) into Equation (50), we have

$$\begin{aligned}
\min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) &= \frac{1}{N} \cdot \sum_{i=1}^N \tilde{l}_f^*(\mathbf{X}_\mu^i) \\
&= \frac{1}{N} \cdot \sum_{i=1}^N -f \left( \frac{dQ}{dP}(\mathbf{X}_\mu^i) \right) \frac{dP}{d\mu}(\mathbf{X}_\mu^i). \tag{55}
\end{aligned}$$

Thus,

$$\begin{aligned}
E_\mu \left[ \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] &= E_\mu \left[ \frac{1}{N} \cdot \sum_{i=1}^N -f \left( \frac{dQ}{dP}(\mathbf{x}_i) \right) \frac{dP}{d\mu}(\mathbf{x}_i) \right] \\
&= -\frac{1}{N} \cdot \sum_{i=1}^N E_\mu \left[ f \left( \frac{dQ}{dP}(\mathbf{x}_i) \right) \frac{dP}{d\mu}(\mathbf{x}_i) \right]
\end{aligned}$$

$$\begin{aligned}
&= -\frac{1}{N} \cdot \sum_{i=1}^N D_f(Q||P) \\
&= -D_f(Q||P),
\end{aligned} \tag{56}$$

From Equations (54) and (56), we have

$$E_\mu \left[ \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi) \right] = \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) = \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu \left[ \mathcal{L}_f^{(R,S)}(\phi) \right], \tag{57}$$

and the equality in each Equation (57) holds for  $\phi(\mathbf{x}) = dQ/dP(\mathbf{x})$ .

This completes the proof.  $\square$

The following theorem presents the convergence rate of the expected value of the distance between two neighboring samples. Similar theorems have been presented in studies on order statistics of multidimensional continuous random variables (e.g., Biau & Devroye (2015), p. 17, Theorem 2.1).

**Theorem C.12** (Theorem 4.3 restated). *Assume that  $\Omega$  is a compact set, as stated in Assumption C.7. Let  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  denote the nearest neighbor of  $\mathbf{x}$  in  $\hat{\mathbf{X}}_{\mu[N]}$ . Specifically, let  $\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  be  $\mathbf{X}_\mu^i$  in  $\hat{\mathbf{X}}_{\mu[N]}$  such that*

$$\|\mathbf{X}_\mu^i - \mathbf{x}\|_\infty < \|\mathbf{X}_\mu^j - \mathbf{x}\|_\infty \ (\forall j < i), \quad \text{and} \quad \|\mathbf{X}_\mu^i - \mathbf{x}\|_\infty \leq \|\mathbf{X}_\mu^j - \mathbf{x}\|_\infty \ (\forall j > i). \tag{58}$$

Additionally, let  $\text{diag}(\Omega)$  denote the diameter of  $\Omega$ . i.e,  $\text{diag}(\mathcal{B}) = \inf_{r \in \mathbb{R}} \{ \mathcal{B} \subseteq \Delta(\mathbf{a}, r) \mid \exists \mathbf{a} \in \mathcal{B} \}$ , where  $\Delta(\mathbf{a}, r)$  denotes the  $d$ -dimensional interval centered at  $\mathbf{a}$  with each side of length  $r$ :  $\Delta(\mathbf{a}, r) = \{ \mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \mathbf{a}\|_\infty < r/2 \}$ .

Then, for  $1 \leq \kappa \leq d$ ,

$$E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa \leq \text{diag}(\Omega)^\kappa \cdot \left( \frac{1}{N+1} \right)^{\kappa/d}, \quad \text{for all } N \geq 1. \tag{59}$$

*proof of Theorem C.12.* Let we rewrite  $\mathbf{X}$  in Equation (59) as  $\mathbf{X}_\mu^{N+1}$ . Subsequently, let  $\hat{\mathbf{X}}_{\mu[N+1]} = \hat{\mathbf{X}}_{\mu[N]} \cup \{ \mathbf{X}_\mu^{N+1} \}$ . Let  $\Delta_i = \Omega \cap \Delta(\mathbf{X}_\mu^i, \|\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i\|_\infty)$ , where  $\Delta(\mathbf{a}, r) = \{ \mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x} - \mathbf{a}\|_\infty < r/2 \}$ . Note that,  $\Delta_i \cap \Delta_j = \emptyset$  if  $i \neq j$ . Thus,  $\sqcup_{i=1}^{N+1} \Delta_i \subseteq \Omega$ .

Now, let  $\lambda$  denote the Lebesgue measure on  $\mathbb{R}^d$ . Then, we have

$$\sum_{i=1}^{N+1} \lambda(\Delta_i) = \lambda(\sqcup_{i=1}^{N+1} \Delta_i) \leq \lambda(\Omega) \leq \text{diag}(\Omega)^d, \tag{60}$$

Subsequently, since  $\lambda(\Delta_i) = \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^d$ , we have

$$\sum_{i=1}^{N+1} \lambda(\Delta_i) = \sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^d. \tag{61}$$

Thus, from Equations (60) and (61), we have

$$\sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^d \leq \text{diag}(\Omega)^d. \tag{62}$$

Note that, it follows from Jensen's inequality that

$$\frac{1}{N+1} \sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^\kappa \leq \left\{ \frac{1}{N+1} \sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^d \right\}^{\kappa/d}. \tag{63}$$

From Equations (62) and (63), we have

$$\begin{aligned}
 \frac{1}{N+1} \sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^\kappa &\leq \left\{ \frac{1}{N+1} \sum_{i=1}^{N+1} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^d \right\}^{\kappa/d} \\
 &\leq \left\{ \frac{1}{N+1} \cdot \text{diag}(\Omega)^d \right\}^{\kappa/d} \\
 &= \text{diag}(\Omega)^\kappa \cdot \left( \frac{1}{N+1} \right)^{\kappa/d}.
 \end{aligned} \tag{64}$$

Thus,

$$\frac{1}{N+1} \sum_{i=1}^{N+1} E_{\mathbf{X}_\mu^i} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa \leq \text{diag}(\Omega)^\kappa \cdot \left( \frac{1}{N+1} \right)^{\kappa/d}, \tag{65}$$

where  $E_{\mathbf{X}_\mu^i} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa$  denotes the expectation of  $\left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{X}_\mu^i) - \mathbf{X}_\mu^i \right\|_\infty^\kappa$  with respect to  $\mathbf{X}_\mu^i$ .

Note that,

$$E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa = E_{\mathbf{X}_\mu^i} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa. \tag{66}$$

Therefore,

$$E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa = \frac{1}{N+1} \sum_{i=1}^{N+1} E_{\mathbf{X}_\mu^i} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa. \tag{67}$$

Finally, from Equations (65) and (67), we have

$$E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa = \frac{1}{N+1} \sum_{i=1}^{N+1} E_{\mathbf{X}_\mu^i} \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^\kappa \leq \text{diag}(\Omega)^\kappa \cdot \left( \frac{1}{N+1} \right)^{\kappa/d}. \tag{68}$$

This completes the proof.  $\square$

**Corollary C.13.** Assume the same assumption as in Theorem C.12. Then, for  $1 \leq p \leq d$ ,

$$\overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_\mu \left[ \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \leq \text{diag}(\Omega). \tag{69}$$

*proof of Corollary C.13.* First, from Theorem C.12 when  $\kappa = p$ ,

$$E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \leq \text{diag}(\Omega)^p \cdot \left( \frac{1}{N+1} \right)^{p/d}, \quad \text{for all } N \geq 1. \tag{70}$$

Thus, for all  $N \geq 1$ ,

$$\begin{aligned}
 \left\{ E_\mu \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right\}^{1/p} &\leq \left\{ \text{diag}(\Omega)^p \cdot \left( \frac{1}{N+1} \right)^{p/d} \right\}^{1/p} \\
 &= \text{diag}(\Omega) \cdot \left( \frac{1}{N+1} \right)^{1/d}
 \end{aligned} \tag{71}$$

Taking  $\overline{\lim}_{N \rightarrow \infty}$  on both sides of the above inequality, we have

$$\begin{aligned}
 \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_\mu \left[ \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \\
 &\leq \overline{\lim}_{N \rightarrow \infty} \left\{ N^{1/d} \cdot \text{diag}(\Omega) \cdot \left( \frac{1}{N+1} \right)^{1/d} \right\} \\
 &= \text{diag}(\Omega).
 \end{aligned} \tag{72}$$

This completes the proof.  $\square$

1134 **Corollary C.14.** Assume the same assumption as in Theorem C.12. Then, for  $1 \leq p \leq d/2$ ,

$$\begin{aligned} & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\ & \leq \text{diag}(\Omega) \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)}. \end{aligned} \quad (73)$$

1142 proof of Corollary C.14. First from Theorem C.12 when  $\kappa = 2 \cdot p$  and  $\mu = P$ ,

$$E_P \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \leq \text{diag}(\Omega)^{2 \cdot p} \cdot \left( \frac{1}{N+1} \right)^{2 \cdot p / d}, \quad \text{for all } N \geq 1. \quad (74)$$

1147 Thus, for all  $N \geq 1$ ,

$$\begin{aligned} & \left\{ E_P \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \right\}^{1/(2 \cdot p)} \leq \left\{ \text{diag}(\Omega)^{2 \cdot p} \cdot \left( \frac{1}{N+1} \right)^{2 \cdot p / d} \right\}^{1/(2 \cdot p)} \\ & = \text{diag}(\Omega) \cdot \left( \frac{1}{N+1} \right)^{1/d} \end{aligned} \quad (75)$$

1154 Now, using Hölder's inequality, we have

$$\begin{aligned} & E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \\ & \leq \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)} \cdot \left( E_P \left[ \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)} \\ & \leq \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)} \cdot \text{diag}(\Omega) \cdot \left( \frac{1}{N+1} \right)^{1/d} \end{aligned} \quad (76)$$

1165 Taking  $\overline{\lim}_{N \rightarrow \infty}$  on both sides of the above inequality, we have

$$\begin{aligned} & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\ & \leq \overline{\lim}_{N \rightarrow \infty} \left\{ N^{1/d} \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)} \cdot \text{diag}(\Omega) \cdot \left( \frac{1}{N+1} \right)^{1/d} \right\} \\ & = \text{diag}(\Omega) \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right)^{1/(2 \cdot p)} \end{aligned} \quad (77)$$

1177 This completes the proof.  $\square$

1178 **Lemma C.15.** Let  $\mu$  be a probability measure on  $\mathbb{R}^d$  with  $d \geq 1$ . Assume that  $\mu \ll \lambda$ , where  $\lambda$  denotes the Lebesgue measure on  $\mathbb{R}^d$ . Let  $\|\cdot\|_{\infty}$  denote the maximum norm in  $\mathbb{R}^d$ :  $\|\mathbf{y} - \mathbf{x}\|_{\infty} = \max_{1 \leq i \leq d} |y^i - x^i|$ , where  $\mathbf{y} = (y^1, y^2, \dots, y^N)$  and  $\mathbf{x} = (x^1, x^2, \dots, x^N)$ . Additionally, let  $\Delta(\mathbf{x}, r)$  denote the  $d$ -dimensional interval centered at  $\mathbf{x}$  with each side of length  $r$ :  $\Delta(\mathbf{x}, r) = \{\mathbf{x}' \in \mathbb{R}^d \mid \|\mathbf{x}' - \mathbf{x}\|_{\infty} \leq r/2\}$ .

1184 Then, for any interior point  $\mathbf{x}$  in  $\Omega$ ,

$$\mu(\Delta(\mathbf{x}, r)) = \frac{d\mu}{d\lambda}(\mathbf{x}) \cdot r^d + o(r^d), \quad \text{as } r \rightarrow 0, \quad (78)$$

1185 where  $f(r) = o(g(r))$ , as  $r \rightarrow 0$ , denotes asymptotic domination such that  $\lim_{r \rightarrow 0} f(r)/g(r) = 0$ .

proof of Lemma C.15. Note that, if  $\mathbf{x}$  is an interior point in  $\Omega$ , it holds that

$$\lim_{r \rightarrow \infty} \frac{\mu(\Delta(\mathbf{x}, r))}{\lambda(\Delta(\mathbf{x}, r))} = \frac{d\mu}{d\lambda}(\mathbf{x}). \quad (79)$$

From Equation (79), we have

$$\begin{aligned} \lim_{r \rightarrow \infty} \frac{\mu(\Delta(\mathbf{x}, r))}{r^d} &= \lim_{r \rightarrow \infty} \frac{\mu(\Delta(\mathbf{x}, r))}{r^d} \\ &= \lim_{r \rightarrow \infty} \frac{\mu(\Delta(\mathbf{x}, r))}{\lambda(\Delta(\mathbf{x}, r))} \end{aligned} \quad (80)$$

$$= \frac{d\mu}{d\lambda}(\mathbf{x}). \quad (81)$$

Here, we use an equation such that  $\lambda(\Delta(\mathbf{x}, r)) = r^d$  in Equation 80.

From Equation (81), we observe that

$$\mu(\Delta(\mathbf{x}, r)) = \frac{d\mu}{d\lambda}(\mathbf{x}) \cdot r^d + o(r^d), \quad \text{as } r \rightarrow 0. \quad (82)$$

This completes the proof.  $\square$

**Corollary C.16.** Assume the same assumptions as in Lemma C.15. Let  $\mathbf{X}$  be a random variable drawn from  $\mu$ , and let  $E_{\mathbf{X}}$  denote the expectation with respect to  $\mathbf{X}$ .

Then, for any interior point  $\mathbf{x}_0$  in  $\Omega$ ,

$$E_{\mathbf{X}} \left[ \|\mathbf{x}_0 - \mathbf{X}\|_{\infty}^p \cdot I(\Delta(\mathbf{x}_0, r))(\mathbf{X}) \right] = \frac{d\mu}{d\lambda}(\mathbf{x}_0) \cdot r^{p+d+1} + o(r^{p+d+1}), \quad \text{as } r \rightarrow 0, \quad (83)$$

where  $I(A)(\cdot)$  is the indicator function for  $A$ :  $I(A)(\mathbf{x}) = 1$  if  $\mathbf{x} \in A$ , and 0 otherwise.

proof of Corollary C.16. Consider the integration variable from  $\mathbf{x}$  to  $r$  such that

$$\|\mathbf{x}_0 - \mathbf{x}\|_{\infty}^p = r. \quad (84)$$

Then, from Lemma C.15, we have, as  $r \rightarrow 0$ ,

$$I(\Delta(\mathbf{x}_0, r))(\mathbf{x}) \cdot \frac{d\mu}{d\lambda}(\mathbf{x}) d\mathbf{x} = \frac{d\mu}{d\lambda}(\mathbf{x}_0) \cdot r^d + o(r^d). \quad (85)$$

From the definition of expectation with the density  $d\mu/d\lambda$  and Equation (85), we have, as  $r \rightarrow 0$ ,

$$\begin{aligned} E_{\mathbf{X}} \left[ \|\mathbf{x}_0 - \mathbf{X}\|_{\infty}^p \cdot I(\Delta(\mathbf{x}_0, r))(\mathbf{X}) \right] \\ = \int \|\mathbf{x}_0 - \mathbf{x}\|_{\infty}^p \cdot I(\Delta(\mathbf{x}_0, r))(\mathbf{x}) \cdot \frac{d\mu}{d\lambda}(\mathbf{x}) d\mathbf{x} \\ = \int r^p \cdot \left( \frac{d\mu}{d\lambda}(\mathbf{x}_0) \cdot r^d + o(r^d) \right) dr \\ = \frac{d\mu}{d\lambda}(\mathbf{x}_0) \cdot r^{p+d+1} + o(r^{p+d+1}). \end{aligned} \quad (86)$$

This completes the proof.  $\square$

**Theorem C.17** (Theorem 4.4 restated). Let  $P$  and  $Q$  be probability measures on a compact set  $\Omega$  in  $\mathbb{R}^d$  with  $d \geq 1$ . Assume that  $P \ll \lambda$  and  $Q \ll \lambda$ , where  $\lambda$  denotes the Lebesgue measure on  $\mathbb{R}^d$ . Let  $p$  be positive constant such that  $p \geq 1$ . Assume  $E[(dQ/dP)^p] < \infty$ .

Then,

$$\varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{P[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \|\mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x}\|_{\infty}^p \right] \right] \right\}^{1/p}$$

$$\geq e^{-1} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p}, \quad (87)$$

where  $E_{\hat{\mathbf{X}}_{P[N]}}[\cdot]$  denotes the expectation on each variable in  $\hat{\mathbf{X}}_{P[N]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^N\}$ .

*proof of Theorem C.17.* Let

$$B_i = \left\{ \mathbf{x} \in \Omega \mid \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty \leq \left( \frac{1}{N} \right)^{1/d} \right\}, \quad (88)$$

where  $\text{diag}(\Omega)$  denotes the diameter of  $\Omega$ :  $\text{diag}(\mathcal{B}) = \inf_{r \in \mathbb{R}} \{ \mathcal{B} \subseteq \Delta(\mathbf{a}, r) \mid \exists \mathbf{a} \in \mathcal{B} \}$ .

Since  $\mathbf{X}_{P[N]}^{(1)}(\mathbf{x})$  is the nearest neighbor in  $\{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^N\}$  for  $\mathbf{x}$ ,

$$\begin{aligned} 1 \leq \exists i \leq N \text{ s.t. } \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty &\leq \left( \frac{1}{N} \right)^{1/d} \\ \iff \quad \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty &\leq \left( \frac{1}{N} \right)^{1/d} \end{aligned} \quad (89)$$

Thus,

$$\begin{aligned} &\left\{ \mathbf{x} \in \Omega \mid \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty \leq \left( \frac{1}{N} \right)^{1/d} \right\} \\ &= \bigcup_{i=1}^N \left\{ \mathbf{x} \in \Omega \mid \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty \leq \left( \frac{1}{N} \right)^{1/d} \right\} = \bigcup_{i=1}^N B_i \end{aligned} \quad (90)$$

Next, define

$$Z_N(\mathbf{x}) = \sum_{i=1}^N I(B_i)(\mathbf{x}). \quad (91)$$

Let  $\mathbf{X}_P$  be a random variable drawn from  $P$  with  $\mathbf{X}_P \perp\!\!\!\perp \mathbf{X}_P^i$ , for  $1 \leq i \leq N$ .

From Lemma C.15,

$$\begin{aligned} P(I(B_i)(\mathbf{X}_P) = 1) &= P(B_i) \\ &= \frac{dP}{d\lambda}(\mathbf{X}_P) \cdot \left( \frac{1}{N^{1/d}} \right)^d + o\left(\frac{1}{N^{1/d}}\right)^d \\ &= \frac{dP}{d\lambda}(\mathbf{X}_P) \cdot \frac{1}{N} + o\left(\frac{1}{N}\right) \\ &= \frac{1}{N} + o\left(\frac{1}{N}\right), \end{aligned} \quad (92)$$

and  $I(B_i)(\mathbf{X}_P) \in \{0, 1\}$  and  $I(B_i)(\mathbf{X}_P) \perp\!\!\!\perp I(B_j)(\mathbf{X}_P)$  for  $i \neq j$ . Namely,  $Z_N(\mathbf{X}_P)$  follows a binomial distribution with the number of trials  $N$  and success probability for each trial  $1/N$ .

Then, we obtain

$$\begin{aligned} E_{\hat{\mathbf{X}}_{P[N]}}[I(\{Z_N(\mathbf{X}_P) = 0\})] &= \left( 1 - \frac{1}{N} - o\left(\frac{1}{N}\right) \right)^N \\ &= \left( 1 - \frac{1}{N} - o\left(\frac{1}{N}\right) \right)^{N-1} \end{aligned} \quad (93)$$

Additionally, note that

$$Z_N(\mathbf{x}) \geq I\left(\bigcup_{i=1}^N B_i\right)(\mathbf{x}),$$

1296 and  
1297

1298  $Z_N(\mathbf{x}) \geq 1 \implies I\left(\bigcup_{i=1}^N B_i\right)(\mathbf{x}) = 1.$   
1299  
1300

1301 In particular,

1302  $Z_N(\mathbf{x}) = 1 \implies \sum_{i=1}^N I(B_i)(\mathbf{x}) = 1.$   
1303  
1304  
1305

1306 Therefore,

1307  $Z_N(\mathbf{x}) = 1 \iff \sum_{i=1}^N I(B_i)(\mathbf{x}) = 1.$   
1308  
1309  
1310

1311 Now, we obtain

1312 
$$\begin{aligned} & N^{p/d} \cdot E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \\ & \geq N^{p/d} \cdot E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right. \\ & \quad \times I \left( \left\{ \mathbf{x} \in \Omega \mid \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty \leq \left( \frac{1}{N} \right)^{1/d} \right\} \right) \\ & \quad \times I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right) \\ & = N^{p/d} \cdot E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right. \\ & \quad \times I \left( \bigcup_{i=1}^N B_i \right) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right) \\ & = N^{p/d} \cdot E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right. \\ & \quad \times \sum_{i=1}^N I(B_i) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right) \\ & = N^{p/d} \cdot \sum_{i=1}^N E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right. \\ & \quad \times I(B_i) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right) \\ & = N^{p/d} \cdot \sum_{i=1}^N E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty^p \right. \\ & \quad \times I(B_i) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right). \end{aligned} \tag{95}$$

1343 Now, let

1344  $Z_N^{-j}(\mathbf{x}) = \sum_{i \neq j}^N I(B_i)(\mathbf{x}).$   
1345  
1346  
1347

1348 Then,

1349  $I(B_i) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1 \right\} \right) = I(B_i) \cdot I \left( \left\{ \mathbf{x} \in \Omega \mid Z_N^{-i}(\mathbf{x}) = 0 \right\} \right).$

1350 (96)

1351

Additionally, let  $\hat{\mathbf{X}}_{P[N]}^{-i}$  denote the subset of  $\hat{\mathbf{X}}_{P[N]}$  excluding  $\mathbf{X}_P^i$ . i.e.,  $\hat{\mathbf{X}}_{P[N]}^{-i} = \hat{\mathbf{X}}_{P[N]} \setminus \{\mathbf{X}_P^i\}$ . Let  $E_N^{-i}[\cdot]$  denote the expectation over the variables in  $\hat{\mathbf{X}}_{P[N]}^{-i}$ , which is equivalent to  $E_{\hat{\mathbf{X}}_{P[N]}^{-i}}[\cdot]$ .

From Equation (93),

$$E_N^{-i} \left[ I(B_i) \cdot I\left(\left\{ \mathbf{x} \in \Omega \mid Z_N^{-i}(\mathbf{x}) = 0 \right\}\right) \right] = \left(1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right)\right)^{N-2}. \quad (97)$$

1362 From Equations (96) and (97), we have

$$\begin{aligned}
& E_N^{-i} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \cdot I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1\}) \right] \right] \\
&= E_N^{-i} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \cdot I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x})^{-i} = 0\}) \right] \right] \\
&= E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \cdot E_N^{-i} \left[ I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x})^{-i} = 0\}) \right] \right] \\
&= E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \times \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \right] \\
&= \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \times E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \right]. \tag{98}
\end{aligned}$$

1378 From Corollary C.16, we have

$$\begin{aligned}
& E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \right] \\
&= \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\{ \frac{dP}{d\mu}(\mathbf{X}_P^i) \cdot \left( \frac{1}{N^{1/d}} \right)^{p+d+1} + o\left( \left( \frac{1}{N^{1/d}} \right)^{p+d+1} \right) \right\} \\
&= \frac{dP}{d\mu}(\mathbf{X}_P^i) \cdot \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left( \frac{1}{N} \right)^{1+p/d} + o\left( \left( \frac{1}{N} \right)^{1+p/d} \right) \\
&= \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left( \frac{1}{N} \right)^{1+p/d} + o\left( \left( \frac{1}{N} \right)^{1+p/d} \right).
\end{aligned} \tag{99}$$

1391 From Equations (98) and (99), we obtain

$$\begin{aligned}
& E_N^{-i} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_{\infty}^p \times I(B_i) \cdot I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1\}) \right] \right] \\
&= \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \\
&\quad \times \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left( \frac{1}{N} \right)^{1+p/d} + o\left( \left( \frac{1}{N} \right)^{1+p/d} \right). \tag{100}
\end{aligned}$$

1401 From Equations (95) and (100), we obtain, as  $N \rightarrow \infty$ ,

$$N^{p/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ E_P \left[ \left\{ \frac{dQ}{dP} \left( \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) \right) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right]$$

$$\begin{aligned}
& \geq N^{p/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \sum_{i=1}^N E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty^p \right. \right. \\
& \quad \times I(B_i) \cdot I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1\}) \left. \right] \\
& = \sum_{i=1}^N N^{p/d} \cdot E_{\mathbf{X}_P^i} \left[ E_N^{-i} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left\| \mathbf{X}_P^i - \mathbf{x} \right\|_\infty^p \right. \right. \right. \\
& \quad \times I(B_i) \cdot I(\{\mathbf{x} \in \Omega \mid Z_N(\mathbf{x}) = 1\}) \left. \right] \left. \right] \\
& = \sum_{i=1}^N N^{p/d} \cdot E_{\mathbf{X}_P^i} \left[ \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \right. \\
& \quad \times \left\{ \frac{dQ}{dP} (\mathbf{X}_P^i) \right\}^p \cdot \left( \frac{1}{N} \right)^{1+p/d} + o\left(\left(\frac{1}{N}\right)^{1+p/d}\right) \left. \right] \\
& = N \cdot \left\{ \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \right. \\
& \quad \times E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{x}) \right\}^p \right] \cdot \left( \frac{1}{N} \right) + o\left(\frac{1}{N}\right) \left. \right\} \\
& = \left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{x}) \right\}^p \right] + o(1) \right\}.
\end{aligned} \tag{101}$$

As  $N \rightarrow \infty$ , we observe

$$\left( 1 - \frac{1}{N-1} - o\left(\frac{1}{N-1}\right) \right)^{N-2} \longrightarrow e^{-1}. \tag{102}$$

Then, we obtain, from Equation (101)

$$\begin{aligned}
& \lim_{N \rightarrow \infty} N^{p/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{P[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right] \\
& \geq e^{-1} \cdot E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{x}) \right\}^p \right].
\end{aligned} \tag{103}$$

This completes the proof.  $\square$

**Theorem C.18.** Assume that  $f$  satisfies Assumption C.6. For  $\tilde{\mathcal{L}}_f^{(N)}(\phi)$  defined in Definition C.3, let  $\phi_*^{(N)} = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$ .

Then, for any measurable function  $\phi: \Omega \rightarrow \mathbb{R}_{>0}$ , the following equivalence holds:

$$\begin{aligned}
& \phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i) = O_p\left(\frac{1}{\sqrt{N}}\right), \quad \text{for } 1 \leq i \leq N \\
& \iff \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f(\phi) = O_p\left(\frac{1}{\sqrt{N}}\right),
\end{aligned} \tag{104}$$

where  $\{\mathbf{X}_\mu^1, \mathbf{X}_\mu^2, \dots, \mathbf{X}_\mu^N\}$  is defined in Definition C.3, and  $\tilde{\mathcal{L}}_f(\phi)$  is defined in Lemma C.11.

*proof of Theorem C.18.* First, we enumerate several facts used in this proof.

I. From the Central Limit Theorem, we have:

$$\tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) - E_\mu \left[ \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right] = O_p\left(\frac{1}{\sqrt{N}}\right). \tag{105}$$

1458 II. From Proposition C.10, we have, for all  $\mathbf{x} \in \hat{\mathbf{X}}_{\mu[N]}$ :

$$1459 \quad 1460 \quad 1461 \quad \phi_*^{(N)}(\mathbf{x}) = \frac{dQ}{dP}(\mathbf{x}), \quad (106)$$

1462 where  $\hat{\mathbf{X}}_{\mu[N]}$  is defined in Definition C.3.

1463 III. From Equation (106), it follows that:

$$1464 \quad 1465 \quad 1466 \quad \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) = \tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right), \quad (107)$$

1467 and

$$1468 \quad 1469 \quad 1470 \quad E_\mu\left[\tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right)\right] = E_\mu\left[\tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right)\right]. \quad (108)$$

1471 IV. From Lemma C.11, we have:

$$1472 \quad 1473 \quad 1474 \quad \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) = \bar{\mathcal{L}}_f\left(\frac{dQ}{dP}\right) = E_\mu\left[\tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right)\right]. \quad (109)$$

1475 V. From Lemma C.8, for  $\tilde{l}_f(u; \mathbf{x})$  defined in Equation (28), we obtain:

$$1476 \quad 1477 \quad 1478 \quad \frac{d}{du} \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) = 0, \quad (110)$$

1479 and

$$1480 \quad 1481 \quad 1482 \quad \frac{d^2}{du^2} \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) = f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}). \quad (111)$$

1483 VI. From Theorem C.9, we have:

$$1484 \quad 1485 \quad 1486 \quad \tilde{l}_f(u; \mathbf{x}) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right) = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}) \cdot \left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2 \\ 1487 \quad 1488 \quad 1489 \quad + o\left(\left|u - \frac{dQ}{dP}(\mathbf{x})\right|^2\right), \quad (112)$$

1490 where  $f(a) = o(a)$  (as  $a \rightarrow 0$ ) denotes asymptotic domination such that  $\lim_{a \rightarrow 0} f(a)/a = 0$ .

1491 VII. From the assumption that  $E_P[f''(dQ/dP)] < \infty$  and the strong law of large numbers, it  
1492 follows that:

$$1493 \quad 1494 \quad 1495 \quad f''\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i)\right) \cdot \frac{dP}{d\mu}(\mathbf{X}_\mu^i) = O_p(\sqrt{N}). \quad (113)$$

1496 Now, we show the direction “ $\implies$ ” in Equation (104).

1497 Assume that  $\phi(\mathbf{X}_\mu^i) = \phi_*^{(N)}(\mathbf{X}_\mu^i) + O_p(1/\sqrt{N})$  for  $1 \leq i \leq N$ .

1498 From Equations (112) and (113), we have

$$1499 \quad 1500 \quad 1501 \quad \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) - \tilde{l}_f\left(\phi_*^{(N)}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \\ 1502 \quad 1503 \quad 1504 \quad = \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \\ 1505 \quad 1506 \quad 1507 \quad = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i)\right) \cdot \frac{dP}{d\mu}(\mathbf{X}_\mu^i) \cdot \left|\phi(\mathbf{X}_\mu^i) - \frac{dQ}{dP}(\mathbf{X}_\mu^i)\right|^2 + o\left(\left|\phi(\mathbf{X}_\mu^i) - \frac{dQ}{dP}(\mathbf{X}_\mu^i)\right|^2\right) \\ 1508 \quad 1509 \quad 1510 \quad = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i)\right) \cdot \frac{dP}{d\mu}(\mathbf{X}_\mu^i) \cdot \left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i)\right|^2 + o\left(\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i)\right|^2\right)$$

$$\begin{aligned}
&= O_p \left( \sqrt{N} \right) \cdot O_p \left( \left\{ \frac{1}{\sqrt{N}} \right\}^2 \right) \\
&= O_p \left( \frac{1}{\sqrt{N}} \right).
\end{aligned} \tag{114}$$

Thus, we have:

$$\begin{aligned}
\tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) &= \frac{1}{N} \cdot \sum_{i=1}^N \left\{ \tilde{l}_f(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \right\} \\
&= \frac{1}{N} \cdot \sum_{i=1}^N O_p\left(\frac{1}{\sqrt{N}}\right) \\
&= O_p\left(\frac{1}{\sqrt{N}}\right).
\end{aligned} \tag{115}$$

From Equations (105), (107), (109), and (115), we obtain:

$$\begin{aligned}
&\tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} + \left\{ \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \right\} \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} + \left\{ \tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \right\} \quad (\text{by Equation (107)}) \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} + \left\{ \tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right) - E\left[\tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right)\right] \right\} \quad (\text{by Equation (109)}) \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} + \left\{ \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) - E\left[\tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right)\right] \right\} \quad (\text{by Equation (107)}) \\
&= O_p\left(\frac{1}{\sqrt{N}}\right) + O_p\left(\frac{1}{\sqrt{N}}\right) \quad (\text{by Equations (105) and (115)}) \\
&= O_p\left(\frac{1}{\sqrt{N}}\right).
\end{aligned} \tag{116}$$

Thus, we have proved “ $\implies$ ”.

Next, we prove the direction “ $\Leftarrow$ ” in Equation (104).

Suppose

$$\tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) = O_p\left(\frac{1}{\sqrt{N}}\right). \tag{117}$$

From Equations (105), (109), (108), and (117), we obtain

$$\begin{aligned}
&\tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \right\} + \left\{ \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \right\} + \left\{ E\left[\tilde{\mathcal{L}}_f^{(N)}\left(\frac{dQ}{dP}\right)\right] - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} \quad (\text{by Equation (109)}) \\
&= \left\{ \tilde{\mathcal{L}}_f^{(N)}(\phi) - \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \bar{\mathcal{L}}_f(\phi) \right\} + \left\{ E\left[\tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right)\right] - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) \right\} \quad (\text{by Equation (108)}) \\
&= O_p\left(\frac{1}{\sqrt{N}}\right) + O_p\left(\frac{1}{\sqrt{N}}\right) \quad (\text{by Equations (105) and (117)}) \\
&= O_p\left(\frac{1}{\sqrt{N}}\right).
\end{aligned} \tag{118}$$

From Equation (106), we have

$$\begin{aligned} \tilde{\mathcal{L}}_f^{(N)}(\phi) - \tilde{\mathcal{L}}_f^{(N)}\left(\phi_*^{(N)}\right) &= \frac{1}{N} \sum_{i=1}^N \tilde{l}_f\left(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) - \frac{1}{N} \sum_{i=1}^N \tilde{l}_f\left(\phi_*^{(N)}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \\ &= \frac{1}{N} \sum_{i=1}^N \left\{ \tilde{l}_f\left(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) - \tilde{l}_f\left(\phi_*^{(N)}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \right\}. \end{aligned} \quad (119)$$

From Equations (118) and (119), we have

$$\frac{1}{N} \sum_{i=1}^N \left\{ \tilde{l}_f\left(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \right\} = O_p\left(\frac{1}{\sqrt{N}}\right). \quad (120)$$

Let  $a_N^i = E_P\left[\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(k)}(\mathbf{X}_\mu^i)\right|\right]$ . Since  $\mathbf{X}_\mu^i$  is identically distributed for  $1 \leq i \leq N$ , we have  $a_N^i = a_N^1$  for any  $1 \leq i \leq N$ . Thus, define  $A_N = \sup_{k \geq N} a_k^i = \sup_{k \geq N} a_k^1$ .

Using Chebyshev's inequality, we have for any  $\varepsilon > 0$ ,

$$\begin{aligned} P\left(\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(k)}(\mathbf{X}_\mu^i)\right| / A_N > \frac{1}{\varepsilon}\right) &\leq \frac{\varepsilon \cdot E_P\left[\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(k)}(\mathbf{X}_\mu^i)\right|\right]}{A_N} \\ &\leq \frac{\varepsilon \cdot a_N^i}{A_N} \\ &\leq \varepsilon. \end{aligned} \quad (121)$$

Thus,  $\phi(\mathbf{X}_\mu^i) - \phi_*^{(k)}(\mathbf{X}_\mu^i) = O_p(A_N)$ .

Now, we calculate

$$\begin{aligned} \frac{1}{N} \sum_{i=1}^N \left\{ \tilde{l}_f\left(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) - \tilde{l}_f\left(\phi_*^{(N)}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \right\} &= \frac{1}{N} \sum_{i=1}^N \left\{ \tilde{l}_f\left(\phi(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{X}_\mu^i); \mathbf{X}_\mu^i\right) \right\} \\ &= \frac{1}{N} \sum_{i=1}^N \left\{ \frac{1}{2} \cdot \lambda(\mathbf{X}_\mu^i) \cdot O_p\left(\left|\phi(\mathbf{X}_\mu^i) - \frac{dQ}{dP}(\mathbf{X}_\mu^i)\right|^2\right) + o_p\left(\left|\phi(\mathbf{X}_\mu^i) - \frac{dQ}{dP}(\mathbf{X}_\mu^i)\right|^4\right) \right\} \\ &= \frac{1}{N} \sum_{i=1}^N \left\{ \frac{1}{2} \cdot \lambda(\mathbf{X}_\mu^i) \cdot O_p\left(\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i)\right|^2\right) + o_p\left(\left|\phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i)\right|^4\right) \right\} \\ &= \frac{1}{N} \cdot N \cdot \frac{1}{2} \cdot O_p\left(\sqrt{N}\right) \cdot O_p(A_N^2) + \frac{1}{N} \cdot N \cdot \frac{1}{2} \cdot o_p(A_N^4) \\ &= O_p\left(\sqrt{N}\right) \cdot O_p(A_N^2) + o_p(A_N^4). \end{aligned} \quad (122)$$

Here,  $\mathbf{X} = o_p(a_N)$  denotes the convergence in probability with rate  $a_N$  in  $\mu$  as  $N \rightarrow \infty$ :  $\mathbf{X} = o_p(a_N)$  (as  $N \rightarrow \infty$ )  $\Leftrightarrow \forall \varepsilon, \forall \delta > 0, \exists N(\varepsilon, \delta) > 0$  such that  $\mu(|\mathbf{X}|/a_N \geq \delta) < \varepsilon$  for  $\forall N \geq N(\varepsilon, \delta)$ .

From Equations (120) and (122), we have

$$O_p\left(\frac{1}{\sqrt{N}}\right) \geq O_p\left(\sqrt{N}\right) \cdot O_p(A_N^2) + o_p(A_N^4). \quad (123)$$

From the definition of  $A_N$ , we observe that  $A_N$  decreases as  $N$  increases. Thus,  $\lim_{N \rightarrow \infty} A_N$  exists and  $0 \leq \lim_{N \rightarrow \infty} A_N < \infty$ .

Suppose that  $\lim_{N \rightarrow \infty} A_N > 0$ . Then, we have

$$O_p(\sqrt{N}) \cdot O_p(A_N^2) + o_p(A_N^4) = O_p(\sqrt{N}) + o_p(1). \quad (124)$$

This contradicts Equation (123). Therefore,  $\lim_{N \rightarrow \infty} A_N = 0$ .

From Equation (123), we have

$$\begin{aligned} O_p\left(\frac{1}{N}\right) &\geq O_p(A_N^2) + o_p\left(\frac{A_N^4}{\sqrt{N}}\right) \\ &= O_p(A_N^2). \end{aligned} \quad (125)$$

Thus,  $A_N = O(1/\sqrt{N})$ .

Finally, we have

$$\phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i) = O_p(A_N) = O_p\left(\frac{1}{\sqrt{N}}\right). \quad (126)$$

Here, we have proved the direction “ $\Leftarrow$ ”.

This completes the proof.  $\square$

**Corollary C.19** (Theorem 4.7 restated). *Assume the same assumption as in Theorem C.18. let  $\phi_*^{(N)} = \arg \min_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} \tilde{\mathcal{L}}_f^{(N)}(\phi)$ .*

*Then, for any measurable function  $\phi : \Omega \rightarrow \mathbb{R}_{>0}$ ,*

$$\begin{aligned} \phi(\mathbf{X}_\mu^i) - \phi_*^{(N)}(\mathbf{X}_\mu^i) &= O_p\left(\frac{1}{\sqrt{N}}\right), \quad \text{for } 1 \leq i \leq N. \\ \iff \mathcal{L}_f^{(R,S)}(\phi) - \inf_{\phi: \Omega \rightarrow \mathbb{R}_{>0}} E_\mu [\mathcal{L}_f^{(R,S)}(\phi)] &= O_p\left(\frac{1}{\sqrt{N}}\right), \end{aligned} \quad (127)$$

*where  $\{\mathbf{X}_\mu^1, \mathbf{X}_\mu^2, \dots, \mathbf{X}_\mu^N\}$  is defined in Definition C.3, and  $\mathcal{L}_f^{(R,S)}(\phi)$  is defined in Definition C.2.*

*proof of Corollary C.19.* From Lemma C.11, we have  $\mathcal{L}_f^{(R,S)}(\phi) = \bar{\mathcal{L}}_f(\phi)$ .

Therefore, Equation (127) follows directly from Equation (104).

This completes the proof.  $\square$

**Theorem C.20** (Theorem 4.5 restated). *Assume that  $\Omega$  is a compact set in  $\mathbb{R}^d$  with  $d \geq 3$  and that  $f$  satisfies Assumption C.6. Let  $P$  and  $Q$  be probability measures on  $\Omega$ . Assume that  $P \ll \lambda$  and  $Q \ll \lambda$ , where  $\lambda$  denotes the Lebesgue measure on  $\mathbb{R}^d$ . Let  $T^*(\mathbf{x})$  be the energy function of  $dQ/dP(\mathbf{x})$  defined as  $T^*(\mathbf{x}) = -\log dQ/dP(\mathbf{x})$ .*

*Let  $\tilde{\mathcal{F}}_{K-Lip}^{(N)}$  denote the set of all  $K$ -Lipschitz continuous functions on  $\Omega$  that minimize  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ . Specifically, define*

$$\tilde{\mathcal{F}}^{(N)} = \left\{ \phi_* : \Omega \rightarrow \mathbb{R}_{>0} \mid \tilde{\mathcal{L}}_f^{(N)}(\phi_*) = \min_{\phi} \tilde{\mathcal{L}}_f^{(N)}(\phi) \right\}, \quad (128)$$

*and*

$$\mathcal{F}_{K-Lip} = \left\{ \phi : \Omega \rightarrow \mathbb{R}_{>0} \mid |\phi(\mathbf{y}) - \phi(\mathbf{x})| \leq K \cdot \|\mathbf{y} - \mathbf{x}\|_\infty \text{ for all } \mathbf{y}, \mathbf{x} \in \Omega \right\}. \quad (129)$$

*Subsequently, let*

$$\tilde{\mathcal{F}}_{K-Lip}^{(N)} = \tilde{\mathcal{F}}^{(N)} \cap \mathcal{F}_{K-Lip}. \quad (130)$$

1674 **(Upper Bound)** Assume Assumption C.4: there exists  $L > 0$  such that  $|T^*(\mathbf{y}) - T^*(\mathbf{x})| \leq L \cdot \|\mathbf{y} - \mathbf{x}\|_\infty$  for any  $\mathbf{y}, \mathbf{x} \in \Omega$ , i.e.,  $T^*(\mathbf{x})$  is  $L$ -Lipschitz continuous on  $\Omega$ .

1677 Then, Equation (131) holds for  $1 \leq p \leq d/2$ , such that for any  $\phi \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)}$ ,

$$\begin{aligned} 1679 \quad & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\ 1680 \quad & \leq L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} + K \cdot \text{diag}(\Omega). \end{aligned} \quad (131)$$

1685 **(Lower Bound)** Assume Assumption C.5: there exists  $L > 1$  such that  $(1/L) \cdot \|\mathbf{y} - \mathbf{x}\|_\infty \leq |T^*(\mathbf{y}) - T^*(\mathbf{x})| \leq L \cdot \|\mathbf{y} - \mathbf{x}\|_\infty$  for any  $\mathbf{y}, \mathbf{x} \in \Omega$ , i.e.,  $T^*(\mathbf{x})$  is  $L$ -bi-Lipschitz continuous on  $\Omega$ ; and  $E_P[dQ/dP] < \infty$  with  $1 \leq p \leq d$ .

1688 Then, Equation (132) holds for any  $\phi \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)}$ , such that

$$\begin{aligned} 1690 \quad & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\tilde{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\ 1691 \quad & \geq \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - K \cdot \text{diag}(\Omega) \end{aligned} \quad (132)$$

$$\begin{aligned} 1696 \quad & \geq \frac{1}{L} \cdot e^{\frac{p-1}{p} \cdot KL(Q||P)-1} - K \cdot \text{diag}(\Omega) \end{aligned} \quad (133)$$

1698 proof of Theorem C.20. First, we list the equations used in this proof.

1700 I. By Taylor's theorem for the second-order Taylor polynomial of  $e^{-t}$ , we have

$$e^{-t} = 1 - t + \frac{1}{2} \cdot e^{-c(t)} \cdot t^2, \quad \text{where } 0 \leq |c(t)| \leq |t|. \quad (134)$$

1704 II. From Equation (134), it follows that

$$\begin{aligned} 1706 \quad & \left| \frac{dQ}{dP}(\mathbf{y}) - \frac{dQ}{dP}(\mathbf{x}) \right| \\ 1707 \quad & = e^{-T^*(\mathbf{y})} \cdot \left| 1 - e^{T^*(\mathbf{y}) - T^*(\mathbf{x})} \right| \\ 1708 \quad & = e^{-T^*(\mathbf{y})} \left\{ (T^*(\mathbf{y}) - T^*(\mathbf{x})) + \frac{1}{2} \cdot e^{C(\mathbf{y}, \mathbf{x}, T^*)} \cdot (T^*(\mathbf{y}) - T^*(\mathbf{x}))^2 \right\} \\ 1709 \quad & = \frac{dQ}{dP}(\mathbf{y}) \left\{ (T^*(\mathbf{y}) - T^*(\mathbf{x})) + \frac{1}{2} \cdot e^{C(\mathbf{y}, \mathbf{x}, T^*)} \cdot (T^*(\mathbf{y}) - T^*(\mathbf{x}))^2 \right\}, \\ 1710 \quad & \text{where } 0 \leq |C(\mathbf{y}, \mathbf{x}, T^*)| \leq |T^*(\mathbf{y}) - T^*(\mathbf{x})|. \end{aligned} \quad (135)$$

1716 III. From Corollary C.13, for  $0 \leq p \leq d/2$ ,

$$\begin{aligned} 1718 \quad & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right\}^{1/p} \leq \text{diag}(\Omega). \end{aligned} \quad (136)$$

1722 IV. From Corollary C.14, for  $0 \leq p \leq d/2$ ,

$$\begin{aligned} 1724 \quad & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \\ 1725 \quad & \leq \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} \end{aligned} \quad (137)$$

1728 V. From Equation (137), for  $0 \leq p \leq d/2$ ,

$$\begin{aligned}
& \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \right] \right\}^{1/p} \\
& \leq \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} \left\{ E_P \left[ \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{4 \cdot p} \right] \right\}^{1/(2 \cdot p)} \\
& \leq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} \cdot \text{diag}(\Omega) \cdot \overline{\lim}_{N \rightarrow \infty} \frac{N^{1/d}}{N^{2/d}} \\
& = 0.
\end{aligned} \tag{138}$$

1740 VI. From Theorem C.17, for  $0 \leq p \leq d$ ,

$$\begin{aligned}
& \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{X}_{P[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{P[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right] \right\}^{1/p} \\
& \geq e^{-1} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p},
\end{aligned} \tag{139}$$

1747 where  $E_{\hat{\mathbf{X}}_{P[N]}}[\cdot]$  denotes the expectation on each variable in  $\hat{\mathbf{X}}_{P[N]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^N\}$ .

1750 VII. Let  $\hat{\mathbf{X}}_{\mu[N]}$  denote the set of random variables defined in Proposition C.10. From Proposi-  
1751 tion C.10,

$$\phi \in \tilde{\mathcal{F}}_{K\text{-Lip}}^{(N)} \iff \phi(\mathbf{X}_{\mu}^i) = \frac{dQ}{dP}(\mathbf{X}_{\mu}^i), \quad \text{for } 1 \leq \forall i \leq N. \tag{140}$$

1755 Now, we prove Equation (131). Let  $\phi(\mathbf{x})$  be a member of  $\tilde{\mathcal{F}}_{K\text{-Lip}}^{(N)}$ .

1757 By applying the triangle inequality in the  $L_p$  norm, we have

$$\begin{aligned}
& \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\
& \leq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right\}^{1/p} + \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi(\mathbf{x}) \right|^p \right\}^{1/p}.
\end{aligned} \tag{141}$$

1765 From the  $K$ -Lipschitz continuity of  $\phi$  and Equation (140),

$$\begin{aligned}
& \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} = \left\{ E_P \left| \phi(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \quad (\text{by Equation 140}) \\
& \leq K \cdot \left\{ E_P \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right\}^{1/p}.
\end{aligned} \tag{142}$$

1773 From Equations (136) and (142),

$$\overline{\lim}_{N \rightarrow \infty} \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \leq K \cdot \text{diag}(\Omega). \tag{143}$$

1778 Next, by substituting  $\mathbf{y} = \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  and multiplying by  $\frac{dP}{d\mu}(\mathbf{x})$  in Equation (135), and using the  
1779  $L$ -Lipschitz continuity of  $T^*$ , we have

$$\left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \frac{dQ}{dP}(\mathbf{x}) \right|^p \right\}^{1/p}$$

$$\begin{aligned}
&= \left[ E_P \left| \frac{dQ}{dP}(\mathbf{x}) \times \left\{ \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right) \right. \right. \right. \\
&\quad \left. \left. \left. + \frac{1}{2} \cdot e^{C_1(\mathbf{x})} \cdot \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right)^2 \right\} \right|^p \right]^{1/p}, \\
&\quad \text{where } 0 \leq C_1(\mathbf{x}) \leq \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|. \\
&= \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) \times \left\{ \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right) \right\} \right. \right. \\
&\quad \left. \left. + \frac{dQ}{dP}(\mathbf{x}) \times \left\{ \frac{1}{2} \cdot e^{C_1(\mathbf{x})} \cdot \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right)^2 \right\} \right|^p \right\}^{1/p} \\
&\leq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^p \right] \right\}^{1/p} \\
&\quad + \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \frac{1}{2^p} \cdot e^{p \cdot C_1(\mathbf{x})} \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^{2 \cdot p} \right] \right\}^{1/p} \\
&\leq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^p \right] \right\}^{1/p} \\
&\quad + \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right. \right. \\
&\quad \times \frac{1}{2^p} \cdot e^{p \cdot |T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x})|} \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^{2 \cdot p} \left. \right] \right\}^{1/p} \\
&\quad \left( \because C_1(\mathbf{x}) \leq \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right| \right) \\
&\leq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\
&\quad + \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right. \right. \\
&\quad \times \frac{1}{2^p} \cdot e^{p \cdot L \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}} \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \left. \right] \right\}^{1/p} \\
&\leq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\
&\quad + \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right. \right. \\
&\quad \times \frac{1}{2^p} \cdot e^{p \cdot L \cdot \text{diag}(\Omega)} \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \left. \right] \right\}^{1/p} \\
&= L \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^p \right] \right\}^{1/p} \\
&\quad + \frac{1}{2} \cdot e^{L \cdot \text{diag}(\Omega)} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_{\infty}^{2 \cdot p} \right] \right\}^{1/p}
\end{aligned} \tag{144}$$

1834

1835

From Equations (137), (138) and (144), we have

$$\begin{aligned}
& \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right\}^{1/p} \\
& \leq \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot L \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \\
& \quad + \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \frac{1}{2} \cdot e^{L \cdot \text{diag}(\Omega)} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \right] \right\}^{1/p} \\
& = L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)}. \tag{145}
\end{aligned}$$

Finally, from Equations (143), (141), and (145), we have

$$\begin{aligned}
& \lim_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\
& \leq L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} + \text{diag}(\Omega) \cdot K. \tag{146}
\end{aligned}$$

Thus, it is shown that Equation (131) holds.

Next, we prove Equation (132). By applying the triangle inequality in the  $L_p$  norm, we have

$$\begin{aligned}
& \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\
& \geq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right\}^{1/p} - \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \phi(\mathbf{x}) \right|^p \right\}^{1/p}. \tag{147}
\end{aligned}$$

By substituting  $\mathbf{y} = \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})$  and multiplying by  $\frac{dP}{d\mu}(\mathbf{x})$  in Equation (135) and the  $L$ -bi-Lipschitz continuity of  $T^*$ , we have

$$\begin{aligned}
& \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - \frac{dQ}{dP}(\mathbf{x}) \right|^p \right\}^{1/p} \\
& = \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right. \right. \\
& \quad \times \left. \left. \left\{ \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right) \right. \right. \right. \\
& \quad \left. \left. \left. + \frac{1}{2} \cdot e^{C_1(\mathbf{x})} \cdot \left( T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right)^2 \right\}^p \right\}^{1/p} \right. \\
& \quad \text{where } 0 \leq C_1(\mathbf{x}) \leq \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right| \\
& \geq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^p \right] \right\}^{1/p} \\
& \quad - \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \right. \right. \\
& \quad \left. \left. \times \frac{1}{2^p} \cdot e^{p \cdot |T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x})|} \cdot \left| T^*(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) - T^*(\mathbf{x}) \right|^{2 \cdot p} \right] \right\}^{1/p} \\
& \geq \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \frac{1}{L^p} \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p}
\end{aligned}$$

$$\begin{aligned}
& - \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \right. \right. \\
& \quad \times \frac{1}{2^p} \cdot e^{p \cdot L \cdot \|\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x}\|_\infty} \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \left. \right] \left. \right\}^{1/p} \\
& \geq \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \frac{1}{L^p} \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \\
& \quad - \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \right. \right. \\
& \quad \times \frac{1}{2^p} \cdot e^{p \cdot L \cdot \text{diag}(\Omega)} \cdot L^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \left. \right] \left. \right\}^{1/p} \\
& = \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right\}^{1/p} \\
& \quad - \frac{1}{2} \cdot e^{\text{diag}(\Omega)} \cdot L \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \right] \right\}^{1/p} \tag{148}
\end{aligned}$$

From Equations (137), (138) and (148), we have

$$\begin{aligned}
& \varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left( E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right|^p \right)^{1/p} \right] \right\} \\
& \geq \varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ \frac{1}{L^p} \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right)^{1/p} \right. \right. \\
& \quad \left. \left. - \frac{1}{2} \cdot e^{\text{diag}(\Omega)} \cdot L \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \right] \right)^{1/p} \right] \right\} \\
& \geq \varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ \frac{1}{L} \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right)^{1/p} \right] \right\} \\
& \quad - \varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_{\hat{\mathbf{X}}_{P[N]}} \left[ \frac{1}{2} \cdot e^{\text{diag}(\Omega)} \right. \right. \\
& \quad \times L \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \right] \right)^{1/p} \left. \right] \right\} \\
& \geq \varliminf_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \frac{1}{L} \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^p \right] \right)^{1/p} \right] \\
& \quad - E_{\hat{\mathbf{X}}_{P[N]}} \left[ \varliminf_{N \rightarrow \infty} N^{1/d} \cdot \left\{ \frac{1}{2} \cdot e^{\text{diag}(\Omega)} \right. \right. \\
& \quad \times L \cdot \left( E_P \left[ \left\{ \frac{dQ}{dP} (\mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x})) \right\}^p \cdot \left\| \mathbf{X}_{\mu[N]}^{(1)}(\mathbf{x}) - \mathbf{x} \right\|_\infty^{2 \cdot p} \right] \right)^{1/p} \left. \right] \left. \right] \\
& = e^{-1} \cdot \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p}. \tag{149}
\end{aligned}$$

Finally, from Equations (143), (147), and (149), we have

$$\begin{aligned}
& \varliminf_{N \rightarrow \infty} N^{p/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\
& \geq e^{-1} \cdot \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - \text{diag}(\Omega) \cdot K. \tag{150}
\end{aligned}$$

1944 Thus, it is shown that Equation (132) holds.  
 1945

1946 Next, we prove Equation (133).  
 1947

1948 First, we have  
 1949

$$\begin{aligned} \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} &= \left\{ E_P \left[ \frac{dQ}{dP}(\mathbf{x}) \cdot \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{p-1} \right] \right\}^{1/p} \\ &= \left\{ E_Q \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{p-1} \right] \right\}^{1/p} \\ &= E_Q \left[ e^{\frac{p-1}{p} \cdot \log \frac{dQ}{dP}(\mathbf{x})} \right]. \end{aligned} \quad (151)$$

1950 From Jensen's inequality,  
 1951

$$\begin{aligned} E_Q \left[ e^{\frac{p-1}{p} \cdot \log \frac{dQ}{dP}(\mathbf{x})} \right] &\geq e^{E_Q \left[ \frac{p-1}{p} \cdot \log \frac{dQ}{dP}(\mathbf{x}) \right]} \\ &= e^{\frac{p-1}{p} \cdot KL(Q||P)}. \end{aligned} \quad (152)$$

1952 From Equations (150), (151) and (152),  
 1953

$$\begin{aligned} \lim_{N \rightarrow \infty} N^{p/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\ \geq e^{-1} \cdot \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - \text{diag}(\Omega) \cdot K \\ \geq \frac{1}{L} \cdot e^{\frac{p-1}{p} \cdot KL(Q||P) - 1} - \text{diag}(\Omega) \cdot K. \end{aligned} \quad (153)$$

1954 This completes the proof.  $\square$   
 1955

1956 **Theorem C.21** (Theorem 4.8 restated). *Assume the same assumptions and notations as in Theorem  
 1957 C.20. Additionally, define*

$$\mathcal{F}_{K\text{-}Lip}^{(N)} = \left\{ \phi \in \mathcal{F}_{K\text{-}Lip} \mid \exists \phi_* \in \tilde{\mathcal{F}}_{K\text{-}Lip}^{(N)} \text{ such that } \phi = \phi_* + O_p \left( \frac{1}{\sqrt{N}} \right) \right\}. \quad (154)$$

1958 That is,  $\mathcal{F}_{K\text{-}Lip}^{(N)}$  denotes the set of all functions that differ by at most  $O_p(1/\sqrt{N})$  from some functions  
 1959 that minimize  $\tilde{\mathcal{L}}_f^{(N)}(\cdot)$ .  
 1960

1961 Then, the same results as in Theorem C.20 hold for all  $\phi \in \mathcal{F}_{K\text{-}Lip}^{(N)}$ . Specifically:  
 1962

1963 **(Upper Bound)** Under Assumption C.4, Equation (131) holds for  $1 \leq p \leq d/2$  such that for any  
 1964  $\phi \in \mathcal{F}_{K\text{-}Lip}^{(N)}$ ,

$$\begin{aligned} \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\ \leq L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} + K \cdot \text{diag}(\Omega). \end{aligned} \quad (155)$$

1965 **(Lower Bound)** Under Assumption C.5, Equation (132) holds for any  $\phi \in \mathcal{F}_{K\text{-}Lip}^{(N)}$ , such that  
 1966

$$\underline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right]$$

$$\begin{aligned} & \geq \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - K \cdot \text{diag}(\Omega) \end{aligned} \quad (156)$$

$$\begin{aligned} & \geq \frac{1}{L} \cdot e^{\frac{p-1}{p} \cdot KL(Q||P)^{-1}} - K \cdot \text{diag}(\Omega) \end{aligned} \quad (157)$$

*Proof of Theorem C.21.* First, we prove Equation (155).

Let  $\tilde{\phi}$  be a member of  $\mathcal{F}_{K\text{-Lip}}^{(N)}$ . Then, there exists  $\phi \in \mathcal{F}_{K\text{-Lip}}^{(N)}$  such that  $\tilde{\phi} = \phi + O_p(1/\sqrt{N})$ .

Using the triangle inequality in the  $L_p$  norm, we obtain

$$\begin{aligned} & \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \tilde{\phi}(\mathbf{x}) \right|^p \right\}^{1/p} = \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) + O_p \left( \frac{1}{\sqrt{N}} \right) \right|^p \right\}^{1/p} \\ & \leq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} + \left\{ E_P \left| O_p \left( \frac{1}{\sqrt{N}} \right) \right|^p \right\}^{1/p} \\ & = \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} + O \left( \frac{1}{\sqrt{N}} \right). \end{aligned} \quad (158)$$

From Equations (131) and (159), we have

$$\begin{aligned} & \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \tilde{\phi}(\mathbf{x}) \right|^p \right\}^{1/p} \\ & \leq \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} + O \left( \frac{1}{\sqrt{N}} \right) \right] \\ & = \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} + \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot O \left( \frac{1}{\sqrt{N}} \right) \\ & = \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \\ & = L \cdot \text{diag}(\Omega) \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^{2 \cdot p} \right] \right\}^{1/(2 \cdot p)} + K \cdot \text{diag}(\Omega). \end{aligned} \quad (159)$$

Therefore, Equation (155) is proven.

Next, we prove Equation (156).

By applying the triangle inequality in the  $L_p$  norm, we obtain

$$\begin{aligned} & \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \tilde{\phi}(\mathbf{x}) \right|^p \right\}^{1/p} = \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) + O_p \left( \frac{1}{\sqrt{N}} \right) \right|^p \right\}^{1/p} \\ & \geq \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} - \left\{ E_P \left| O_p \left( \frac{1}{\sqrt{N}} \right) \right|^p \right\}^{1/p} \\ & = \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} - O \left( \frac{1}{\sqrt{N}} \right). \end{aligned} \quad (160)$$

In a similar manner to the derivation of Equation (159), we have

$$\begin{aligned} & \underline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \tilde{\phi}(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\ & \geq \underline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} - O \left( \frac{1}{\sqrt{N}} \right) \right] \\ & = \underline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] - \overline{\lim}_{N \rightarrow \infty} N^{1/d} \cdot O \left( \frac{1}{\sqrt{N}} \right) \end{aligned}$$

$$\begin{aligned}
&= \lim_{N \rightarrow \infty} N^{1/d} \cdot E_{\hat{\mathbf{X}}_{P[N]}} \left[ \left\{ E_P \left| \frac{dQ}{dP}(\mathbf{x}) - \phi(\mathbf{x}) \right|^p \right\}^{1/p} \right] \\
&= \frac{1}{L} \cdot \left\{ E_P \left[ \left\{ \frac{dQ}{dP}(\mathbf{x}) \right\}^p \right] \right\}^{1/p} - K \cdot \text{diag}(\Omega).
\end{aligned} \tag{161}$$

Therefore, Equation (156) is proven.

Equation (157) is obtained in the same manner as in the proof of Theorem C.20.

This completes the proof.  $\square$

## D DETAILS OF THE EXPERIMENTS IN SECTION 3

In this section, we provide details on the experiments reported in Section 3. Each dataset, experimental method, experimental result, and the neural network settings used in the experiments are described in separate subsections.

### D.0.1 DATASETS.

In both experiments investigating the relationship between  $L_p$  errors and KL-divergence in the data, and the relationship between  $L_p$  errors and the dimensionality of the data, the datasets were generated from the following distributions: the numerator distribution is a multidimensional multimodal normal distribution, and the denominator distribution is a multidimensional standard normal distribution.

**Denominator Distribution:** The denominator datasets  $\hat{\mathbf{X}}_{P[R]} = \{\mathbf{X}_P^1, \mathbf{X}_P^2, \dots, \mathbf{X}_P^R\}$  were generated from the following  $d$ -dimensional standard normal distribution:

$$\mathbf{X}_P^i \stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I_d), \tag{162}$$

where  $I_d$  denotes the  $d$ -dimensional identity matrix.

**Numerator Distribution:** The numerator datasets  $\hat{\mathbf{X}}_{Q[S]} = \{\mathbf{X}_Q^1, \mathbf{X}_Q^2, \dots, \mathbf{X}_Q^S\}$  were generated from the following  $d$ -dimensional,  $M$ -multimodal normal distribution:

$$\mathbf{X}_Q^i \stackrel{\text{iid}}{\sim} \prod_{m=1}^M \mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)^{Z_m}, \tag{163}$$

where for each mode  $m$ :

- $Z_m \sim \text{Bernoulli}(1/M)$  and  $\sum_{m=1}^M Z_m = 1$ .
- $\mathbf{r}_m \sim \text{Uniform}(\mathbb{S}^{d-1})$ .

Here,  $\text{Bernoulli}(1/M)$  denotes the Bernoulli distribution with parameter  $1/M$ , and  $\text{Uniform}(\mathbb{S}^{d-1})$  denotes the uniform distribution on the  $d$ -dimensional unit surface  $\mathbb{S}^{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ .

In the aforementioned setting when  $M = 1$ , the KL-divergence of the datasets is calculated as:

$$\begin{aligned}
KL(P||Q) &= E_P \left[ \log \left( \frac{dP}{dQ} \right) \right] \\
&= E_{\mathcal{N}(\mathbf{0}, I_d)} \left[ \log \left( \frac{\mathcal{N}(\mathbf{0}, I_d)}{\mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)} \right) \right] \\
&= \frac{1}{2} \cdot \left[ \log \frac{|\Sigma_p|}{|\Sigma_q|} - d + \text{Tr}(\Sigma_p^{-1} \cdot \Sigma_q) + (\mu_p - \mu_q)^T \cdot \Sigma_p^{-1} \cdot (\mu_p - \mu_q) \right] \\
&= \frac{1}{2} \cdot \left[ \log \frac{|I_d|}{|I_d|} - d + \text{Tr}(I_d \cdot I_d) + (\mu \cdot \mathbf{r}_m)^T \cdot I_d \cdot (\mu \cdot \mathbf{r}_m) \right] \\
&= \frac{1}{2} \cdot (0 - d + d + \mu^2 \cdot \mathbf{r}_m^T \cdot \mathbf{r}_m) \\
&= \frac{1}{2} \cdot \mu^2.
\end{aligned} \tag{164}$$

From Equation (164), the KL-divergence of the datasets for  $M > 1$  is calculated as:

$$\begin{aligned}
 KL(P||Q) &= E_P \left[ \log \left( \frac{dP}{dQ} \right) \right] \\
 &= E_{\mathcal{N}(\mathbf{0}, I_d)} E_{Z_m \sim \text{Bernoulli}(1/M)} \left[ \log \left( \frac{\mathcal{N}(\mathbf{0}, I_d)}{\prod_{m=1}^M \mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)^{Z_m}} \right) \right] \\
 &= E_{\mathcal{N}(\mathbf{0}, I_d)} E_{Z_m \sim \text{Bernoulli}(1/M)} \left[ \log \prod_{m=1}^M \left( \frac{\mathcal{N}(\mathbf{0}, I_d)}{\mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)} \right)^{Z_m} \right] \\
 &= E_{\mathcal{N}(\mathbf{0}, I_d)} E_{Z_m \sim \text{Bernoulli}(1/M)} \left[ \sum_{m=1}^M \log \left( \frac{\mathcal{N}(\mathbf{0}, I_d)}{\mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)} \right) \right] \\
 &= E_{\mathcal{N}(\mathbf{0}, I_d)} \left[ \log \left( \frac{\mathcal{N}(\mathbf{0}, I_d)}{\mathcal{N}(\mu \cdot \mathbf{r}_m, I_d)} \right) \right] \\
 &= \frac{1}{2} \cdot \mu^2. \tag{165}
 \end{aligned}$$

Thus, we set  $\mu = \sqrt{2 \cdot KL(P||Q)}$  in Equation (163) for  $M = 1, 2, 3$ , and  $4$ , where  $KL(P||Q)$  denotes the KL-divergence of the datasets.

## D.1 EXPERIMENTAL PROCEDURE.

We trained neural networks using the training datasets by optimizing KL-divergence and  $\alpha$ -divergence loss functions. Details of the two functions used in the experiments are provided below.

**KL-divergence loss function.** We used the following KL-divergence loss function,  $\mathcal{L}_{\text{KL}}(\cdot)$ , in our experiments:

$$\begin{aligned}
 \mathcal{L}_{\text{KL}}(T) &= \hat{E}_P [e^T] - \hat{E}_Q [T] \\
 &= \frac{1}{S} \cdot \sum_{i=1}^S e^{T(\mathbf{x}_Q^i)} - \frac{1}{R} \cdot \sum_{i=1}^R T(\mathbf{x}_P^i). \tag{166}
 \end{aligned}$$

**$\alpha$ -divergence loss function.** We utilized an  $\alpha$ -divergence loss function proposed in a separate unpublished study, currently under anonymous review. The  $\alpha$ -divergence loss function is defined as:

$$\begin{aligned}
 \mathcal{L}_{\alpha\text{-divergence}}^{(R,S)}(T; \alpha) &= \frac{1}{\alpha} \cdot \hat{E}_{Q[S]} \left[ e^{\alpha \cdot T_\theta} \right] + \frac{1}{1-\alpha} \cdot \hat{E}_{P[R]} \left[ e^{(\alpha-1) \cdot T_\theta} \right] \\
 &= \frac{1}{\alpha} \cdot \frac{1}{S} \cdot \sum_{i=1}^S e^{\alpha \cdot T(\mathbf{x}_Q^i)} + \frac{1}{1-\alpha} \cdot \frac{1}{R} \cdot \sum_{i=1}^R e^{(\alpha-1) \cdot T(\mathbf{x}_P^i)}. \tag{167}
 \end{aligned}$$

For further details and theoretical derivations of the loss function, we refer the reader to the anonymized supplementary material included in this submission (see Anonymous (2024)). This material contains a full explanation of the theoretical framework and the optimization process of the loss function used here.

**$L_p$  Errors vs. KL-Divergence in Data.** We initially created 100 training, validation, and test datasets, each consisting of 10000 samples, with a data dimensionality of 5 and KL-divergence values of 1, 2, 4, 8, 10, 12, and 14, and the numerator datasets of modalities of 1, 2, 3, and 4. The numerator datasets had modalities of 1, 2, 3, and 4, generated from the aforementioned distributions. We trained neural networks using the training datasets by optimizing both the  $\alpha$ -divergence and KL-divergence loss functions. Training was halted if the validation loss, measured using the validation datasets, did not improve over an entire epoch. After training the neural networks, we measured the  $L_p$  errors of the estimated density ratios for  $p = 1, 2$ , and  $3$ , using the test datasets. A total of 100 trials were conducted, and we reported the median  $L_p$  errors along with the interquartile range (25th to 75th percentiles) for each KL-divergence and  $\alpha$ -divergence function.

**$L_p$  Errors vs. the Dimensions of Data.** We initially created 100 training datasets, each consisting of 20000 samples, and 100 validation and test datasets, each consisting of 5000 samples, with data dimensionalities of 50, 100, and 200, and a KL-divergence value of 3. We trained neural networks using the training datasets of sizes 1000, 2000, 4000, 8000, and 16000, by optimizing both the  $\alpha$ -divergence and KL-divergence loss functions. The numerator datasets had modalities of 1, 2, 3, and 4, generated from the aforementioned distributions. Training was halted if the validation loss, measured using the validation datasets, did not improve over an entire epoch. After training the neural networks, we measured the  $L_p$  errors of the estimated density ratios for  $p = 1, 2$ , and 3, using the test datasets. A total of 100 trials were conducted, and we reported the median  $L_p$  errors along with the interquartile range (25th to 75th percentiles) for each KL-divergence and  $\alpha$ -divergence function.

## D.2 RESULTS.

**$L_p$  Errors vs. the KL-Divergence in Data.** The results for each multimodal case  $M = 1, 2, 3$ , and 4 of the numerator datasets are shown in Figure 3. The results of  $M = 1$  were reported in Section 3.

As shown in Figure 3, the estimation errors for  $p > 0$  increased significantly, which accelerates as  $p$  becomes larger. In contrast, when  $p = 0$ , a relatively mild increase was observed. As indicated by Theorem 3.5, these results highlight the impact of the KL-divergence in the data on  $L_p$  error with  $p > 1$  in DRE  $f$ -divergence loss functions. Additionally, little difference was observed in the results among the modalities of the numerator datasets.

**$L_p$  Errors vs. the Dimensions of Data.** The results for each multimodal case  $M = 1, 2, 3$ , and 4 of the numerator datasets are shown in Figure 4 and 5. The results of  $M = 1$  (the first and second rows in Figure 4) were reported in Section 3.

As shown in Figure 2, the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE deteriorated as the data dimensionality increases for both the  $\alpha$ -divergence and KL-divergence loss functions. These results indicate that the curse of dimensionality occurs equally across the  $L_p$  errors, as indicated by Theorem 3.5. Additionally, little difference was observed in the results among the modalities of the numerator datasets.

## D.3 NEURAL NETWORK ARCHITECTURE, OPTIMIZATION ALGORITHM, AND HYPERPARAMETERS.

**$L_p$  Errors vs. the KL-Divergence in Data.** The same neural network architecture, optimization algorithm, and hyperparameters were used for both the KL-divergence and  $\alpha$ -divergence loss functions. A 6-layer perceptron with ReLU activation was employed, with each hidden layer consisting of 1024 nodes. For optimization with the both the KL-divergence and  $\alpha$ -divergence loss functions, the learning rate was 0.0001, and the batch size was 128. Early stopping was applied with a patience of 3 epochs, and the maximum number of epochs was set to 5000. the value of  $\alpha$  for the  $\alpha$ -divergence loss function was set to 0.5, Pytorch (Paszke et al., 2017) library in Python was used to implement all models for DRE, with the Adam optimizer (Kingma, 2014) in PyTorch and an NVIDIA T4 GPU used for training the neural networks.

**$L_p$  Errors vs. the Dimensions of Data.** The same neural network architecture, optimization algorithm, and hyperparameters were used for the KL-divergence and  $\alpha$ -divergence loss functions. A 6-layer perceptron with ReLU activation was employed, with each hidden layer consisting of 1024 nodes. For optimization with the both the KL-divergence and  $\alpha$ -divergence loss functions, the learning rate was 0.0001, and the batch size was 128. Early stopping was applied with a patience of 1 epochs, and the maximum number of epochs was set to 5000. the value of  $\alpha$  for the  $\alpha$ -divergence loss function was set to 0.5, Pytorch (Paszke et al., 2017) library in Python was used to implement all models for DRE, with the Adam optimizer (Kingma, 2014) in PyTorch and an NVIDIA T4 GPU used for training the neural networks.

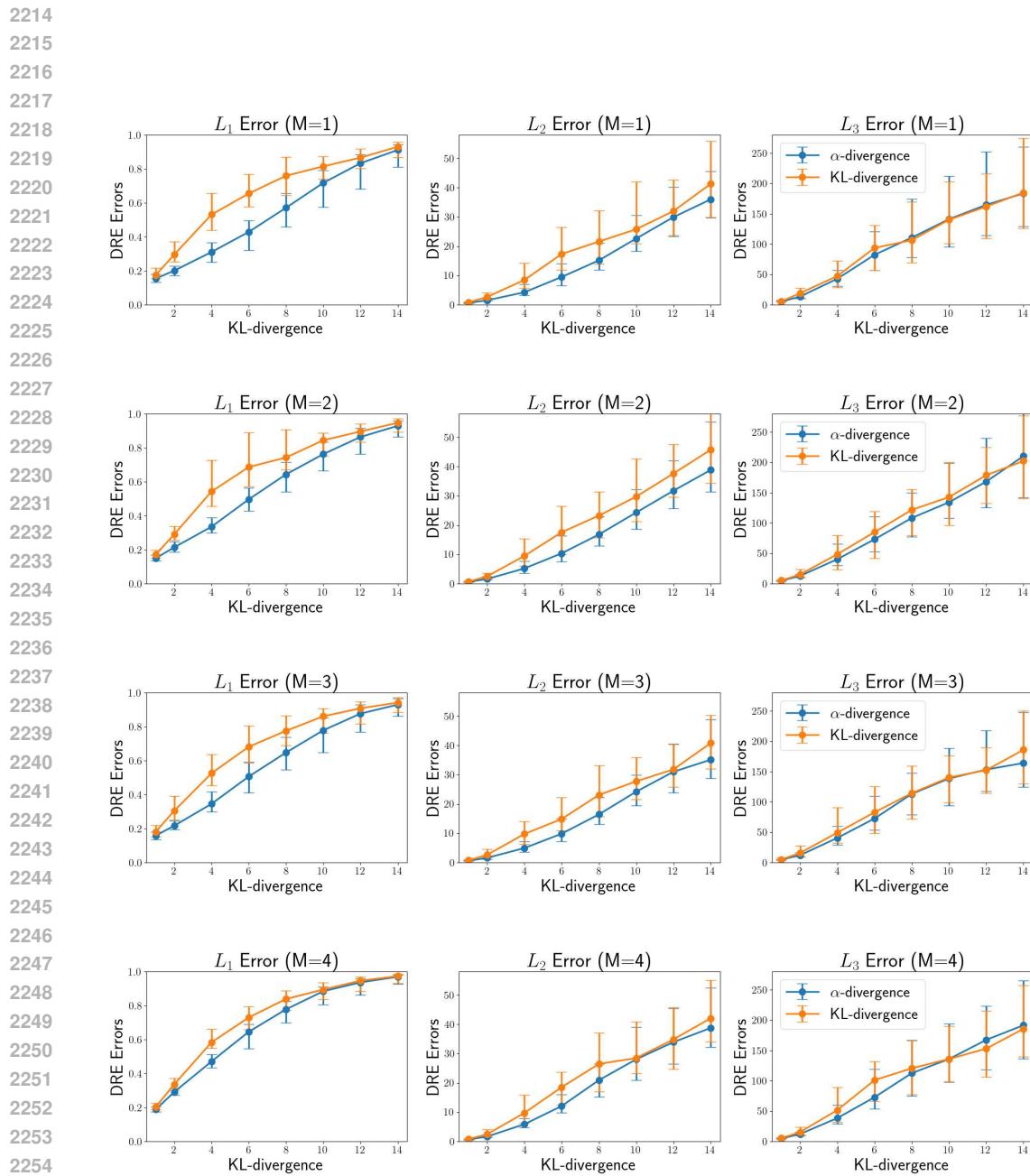


Figure 3: Experimental results of  $L_p$  errors versus the KL-divergence in the data for each multi-modal case  $M = 1, 2, 3$ , and 4 of the numerator datasets, as discussed in Sections 3 and D. The results for  $M = 1$  were reported in Section 3. The  $x$ -axis represents the KL-divergence of synthetic datasets with fixed dimension. The  $y$ -axes of the left, center, and right graphs represent the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE, respectively. The blue line represents errors using the  $\alpha$ -divergence loss function, and the orange line represents errors using the KL-divergence loss function. The error bars represent the interquartile range (25th to 75th percentiles) of the  $y$ -axis values. The plots show the median  $y$ -axis values corresponding to the KL-divergence levels in the synthetic datasets.

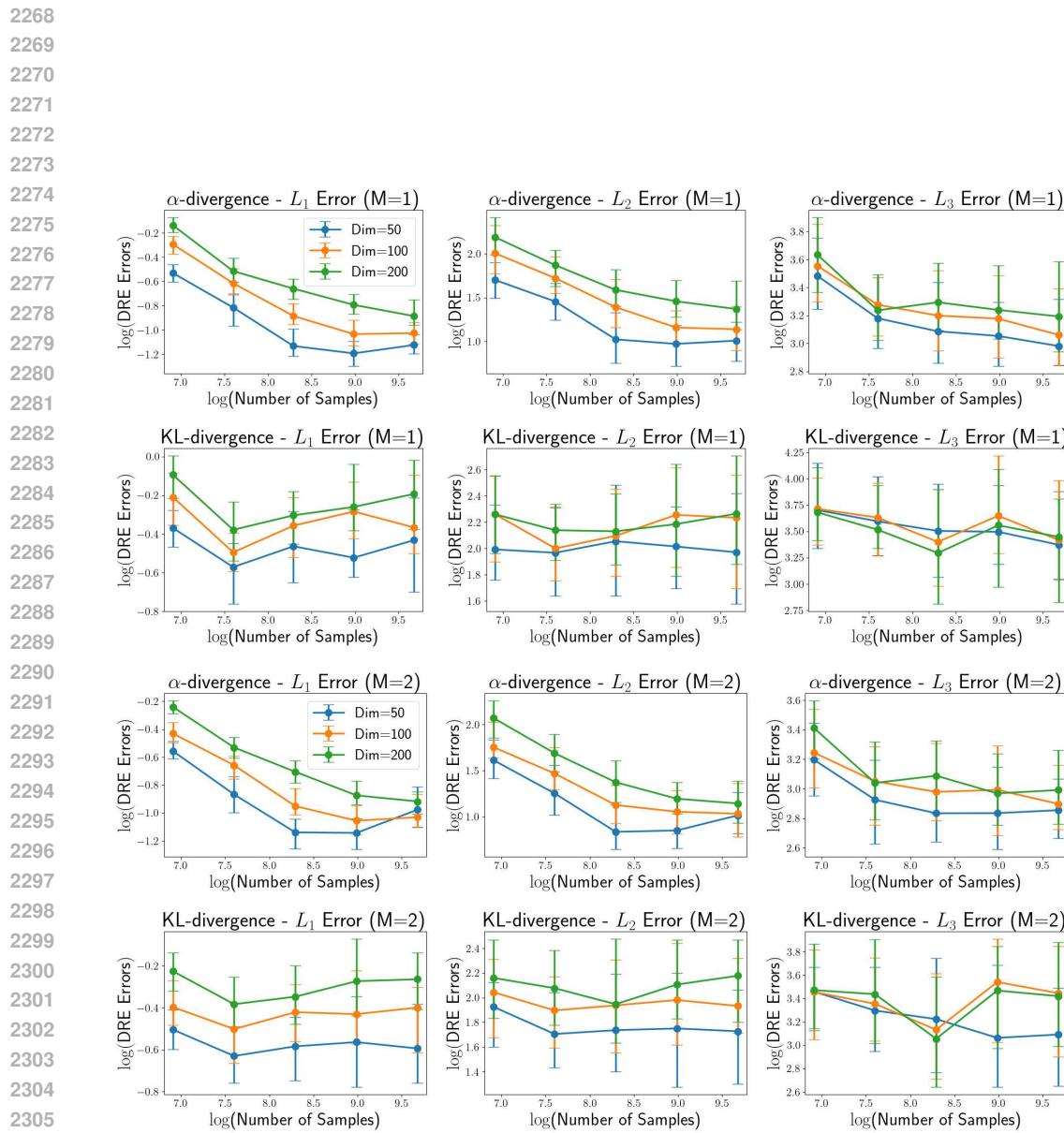


Figure 4: Experimental results of  $L_p$  errors versus the dimensionality of the data for the multimodal case  $M = 1$  and  $2$  in the numerator datasets, as discussed in Sections 3 and D. The results for  $M = 1$  were reported in Section 3. The top row shows the results using the  $\alpha$ -divergence loss function, and the bottom the results using the KL-divergence loss function. The  $x$ -axis represents the logarithm of the number of samples used for the optimizations for DRE. The  $y$ -axes of the left, center, and right graphs represent the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE, respectively. The blue, orange, and green lines represent the results for data dimensionalities of 50, 100, and 200, respectively. The plots show the median  $y$ -axis values, and the error bars indicate the interquartile range (25th to 75th percentiles) of the  $y$ -axis values for the logarithm of the number of samples used for the optimizations of DRE.

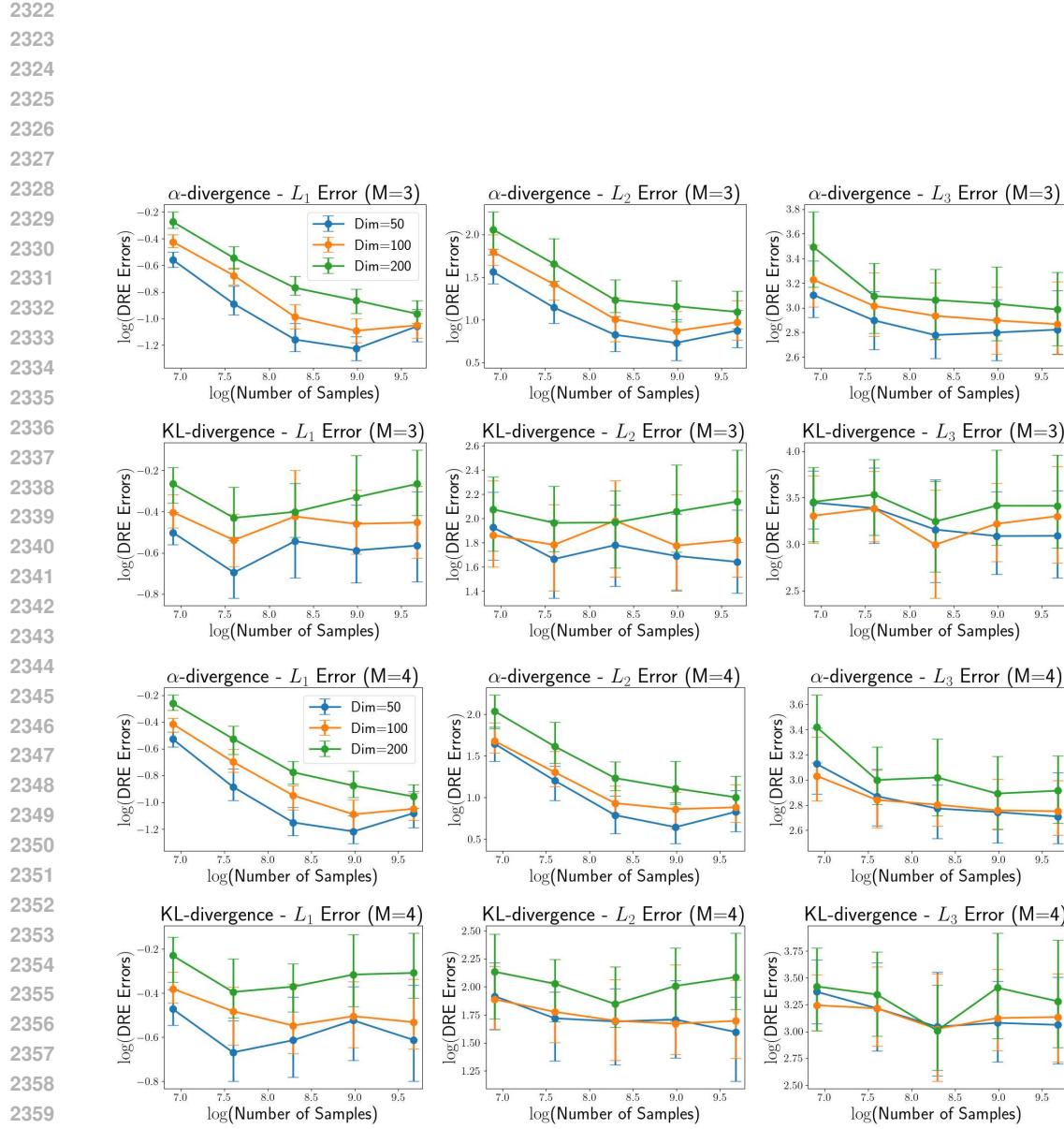


Figure 5: Experimental results of  $L_p$  errors versus the dimensionality of the data for the multimodal case  $M = 3$  and  $4$  in the numerator datasets, as discussed in Sections D. The top row shows the results using the  $\alpha$ -divergence loss function, and the bottom the results using the KL-divergence loss function. The  $x$ -axis represents the logarithm of the number of samples used for the optimizations for DRE. The  $y$ -axes of the left, center, and right graphs represent the  $L_1$ ,  $L_2$ , and  $L_3$  errors in DRE, respectively. Blue, orange, and green lines represent the results for data dimensionalities of 50, 100, and 200, respectively. The plots show the median  $y$ -axis values, and the error bars indicate the interquartile range (25th to 75th percentiles) of the  $y$ -axis values for the logarithm of the number of samples used for the optimizations of DRE.

2376 **E FURTHER DISCUSSIONS RELATED TO THIS STUDY**  
 2377

2378 In this section, we explore further discussions related to this study. First, we compare the upper DRE  
 2379 bound derived in this study with those reported in previous research. Next, we provide remarks on  
 2380 Assumption 3.3, comparing it with related assumptions in prior work. Finally, we highlight the  
 2381 potential applications suggested by this study.  
 2382

2383 **E.1 COMPARISON WITH EXISTING DRE BOUNDS**  
 2384

2385 In this section, we compare our  $L_p$  upper bound in Equation (4) in Theorem 3.5 to known DRE  
 2386 bounds from other methods.  
 2387

2388 The terms related to data dimensionality in our upper bound are tighter than the existing non-  
 2389 parametric minimax upper bounds in DRE. Additionally, to the best of our knowledge, no prior  
 2390 work has provided a term like ours regarding the exponential of the KL-divergence in Equation (6)  
 2391 in Theorem 3.5.  
 2392

2393 Nguyen et al. (2010) presented a minimax upper bound rate of  $O(1/N^{\frac{1}{2+d}})$  for the Hellinger dis-  
 2394 tance between the true and estimated density ratio, obtained by optimizing a KL-divergence loss  
 2395 function. Since the Hellinger distance serves as an upper bound for the total variation distance  
 2396 (Sason & Verdú, 2016), the result from Nguyen et al. (2010) provides an upper bound on the  $L_1$  er-  
 2397 rror in DRE using the KL-divergence loss function. Kanamori et al. (2012) provided an upper bound  
 2398 of  $O(1/N^{\frac{1}{2+d}})$  for DRE using kernel unconstrained least-squares importance fitting (KuLSIF), their  
 2399 proposed DRE method. Under an assumption on the  $\beta$ -Hölder continuity of the probability ratio  
 2400 function, Kpotufe (2017) presented an upper bound of  $O_P(\log N/N^{\frac{\beta}{\beta+d}})$  for DRE using an em-  
 2401 pirical distribution-based estimator, where our case corresponds to  $\beta = 1$ . A recent study (Lin et al.,  
 2402 2023) provided  $L_1$  and  $L_2$  error upper bounds of  $O(1/N^{\frac{1}{2+d}})$  in DRE for an estimator using the  
 2403  $M$ -th nearest neighbor, as  $M$  increases along with the sample size.  
 2404

2405 In terms of comparison with our  $L_p$  lower bound, a minimax  $L_1$  lower bound of  $O(1/N^{\frac{1}{2+d}})$ , for  
 2406 example, was provided by Lin et al. (2023). This lower bound is larger than our lower bound in  
 2407 Equation (5) in Theorem 3.5 and appears tighter than ours. However, minimax lower bounds may  
 2408 not represent the true lower bounds and cannot be directly compared to our lower bound, as discussed  
 2409 in Section 1.  
 2410

2411 **E.2 REMARKS ON ASSUMPTION 3.3 AND RELATED ASSUMPTIONS IN PRIOR WORK**  
 2412

2413 In this section, through a comparison to related assumptions in prior work, we provide remarks on  
 2414 Assumption 3.3 presented in Section 3.  
 2415

2416 We believe that a closely related assumption to Assumption 3.3 can be found in the pseudo self-  
 2417 concordance of losses introduced in Bach (2010). The pseudo self-concordance of losses ensures a  
 2418 deterministic strength of the convexity of loss functions, whereas Assumption 3.3 only provides a  
 2419 statistical measure of the strength of the convexity of the  $f$ -divergence loss functions.  
 2420

2421 Specifically, let  $l(u)$  denote a convex loss function such that  $l''(u) > 0$  for  $u > 0$ . Subsequently,  
 2422 let  $F_u(r) = \{l(u+r) - l(u)\}/l''(u)$ . According to Proposition 1 in Bach (2010), the pseudo  
 2423 self-concordance of the losses implies both  $L(r)$ -smoothness and  $1/L(r)$ -strong convexity of  $F_u(r)$   
 2424 within any interval of a fixed length  $r$ , which is determined independently of  $u$ . Thus, both  $L(r)$ -  
 2425 smoothness and  $1/L(r)$ -strong convexity of  $F_u(r)$  are deterministically guaranteed even when  $u$  is  
 2426 randomly selected.  
 2427

2428 On the other hand, from Theorem C.9 in the appendix, Assumption 3.3 describes a localized property  
 2429 of the strength of the loss convexity such that  
 2430

$$\frac{\tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}) + r; \mathbf{x}\right) - \tilde{l}_f\left(\frac{dQ}{dP}(\mathbf{x}); \mathbf{x}\right)}{|r|^2} = \frac{1}{2} \cdot f''\left(\frac{dQ}{dP}(\mathbf{x})\right) \cdot \frac{dP}{d\mu}(\mathbf{x}) + o(|r|^2). \quad (168)$$

2431 Then, let  $F_{u(\mathbf{x})}(r) = \{\tilde{l}_f(u(\mathbf{x}) + r; \mathbf{x}) - \tilde{l}_f(u(\mathbf{x}); \mathbf{x})\} / \frac{d^2}{du^2} \tilde{l}_f(u(\mathbf{x}); \mathbf{x})$ .  
 2432

According to Equation (168), the  $L$ -smoothness and  $1/L$ -strong convexity of  $F_{u(\mathbf{x})}(r)$  on  $\Delta(\mathbf{x}, r)$  with fixed  $r > 0$  can depend on each point  $\mathbf{x}$ , i.e.,  $L = L(\mathbf{x})$ . Here,  $\Delta(\mathbf{x}, r)$  denotes the  $d$ -dimensional interval centered at  $\mathbf{x}$  with each side of length  $r$ :  $\Delta(\mathbf{x}, r) = \{\mathbf{x}' \in \mathbb{R}^d \mid \|\mathbf{x}' - \mathbf{x}\|_\infty < r/2\}$ .

Considering the expectation with respect to  $\mu$  on both sides of this equation, we have

$$E_\mu \left[ \frac{\tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{x}) + r; \mathbf{x} \right) - \tilde{l}_f \left( \frac{dQ}{dP}(\mathbf{x}); \mathbf{x} \right)}{|r|^2} \right] = \frac{1}{2} \cdot E_P \left[ f'' \left( \frac{dQ}{dP}(\mathbf{x}) \right) \right] + o(|r|^2). \quad (169)$$

This implies that the strength of  $F_{u(\mathbf{x})}(r)$  is statistically represented by  $E_P[f''(dQ/dP)]$ .

Additionally, we note that the expression  $E_P[f''(dQ/dP)]$  resembles the form of the Fisher information when  $f(u) = -\log u$ . Thus, as an alternative perspective on Assumption 3.3, we propose that this assumption establishes an information-theoretic bound for estimation using  $f$ -divergence optimization.

### E.3 APPLICATIONS OF THIS STUDY

In this section, we provide a brief discussion of potential applications highlighted by our findings. The following two key applications can be derived from our results.

**Selecting a benchmark index for evaluating DRE methods.** When evaluating the accuracy of DRE methods using synthetic datasets, the root mean squared error (RMSE) or mean squared error (MSE) is recommended rather than the mean absolute error (MAE). Prior works did not carefully consider the differences in their behavior regarding the KL divergence of the datasets. For example, Kimura & Bondell (2024) used MAE, whereas Kato & Teshima (2021) used MSE.

**Fitting the distribution of base noise for  $f$ -GAN and Normalizing Flow.** Optimization of  $f$ -GANs (Nowozin et al., 2016) could benefit from adjusting the base noise distribution to better match the data. Since the optimization of  $f$ -GANs is equivalent to DRE by optimizing the  $f$ -divergence (Uehara et al., 2016), the accuracy of generative models could be improved by fitting the base parametric models to the data in terms of KL divergence minimization (i.e., likelihood maximization). A similar approach could also be applied to the base models in Normalizing Flow (Papamakarios et al., 2021).

2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

2484  
 2485  
 2486  
 2487  
 2488  
 2489  
 2490  
 2491  
 2492  
 2493  
 2494  
 2495  
 2496  
 2497  
 2498  
 2499  
 2500  
 2501  
 2502  
 2503  
 2504  
 2505  
 2506

2507 Table 2: List of  $f'(\phi)$  and  $f^*(f'(\phi))$  in Equation (1) together with convex functions, as discussed  
 2508 Section 2.2. Part of the list of divergences and their convex functions is based on Nowozin et al.  
 2509 (2016).

2510  
 2511  
 2512  
 2513  
 2514  
 2515  
 2516  
 2517  
 2518  
 2519  
 2520  
 2521  
 2522  
 2523  
 2524  
 2525  
 2526  
 2527  
 2528  
 2529  
 2530  
 2531  
 2532  
 2533  
 2534  
 2535  
 2536  
 2537

Name	convex function $f$	$f'(\phi)$	$f^*(f'(\phi))$
KL	$u \cdot \log u$	$\log(\phi) + 1$	$\phi$
Pearson $\chi^2$	$(u - 1)^2$	$2 \cdot \phi - 2$	$\phi^2 - 2$
Squared Hellinger	$(\sqrt{u} - 1)^2$	$1 - \phi^{-1/2}$	$\phi^{1/2} - 1$
GAN	$u \cdot \log u - (u + 1) \cdot \log(u + 1)$	$-\log(1 + \phi^{-1})$	$\log(1 + \phi)$