# Empirical estimation of entropy functionals with confidence

**Kumar Sricharan**, Department of EECS, University of Michigan
**Raviv Raich**, School of EECS, Oregon State University
**Alfred O. Hero III**, Department of EECS, University of Michigan

February 25, 2011

### Abstract

This paper introduces a class of k-nearest neighbor ($k$-NN) estimators called bi-partite plug-in (BPI) estimators for estimating integrals of non-linear functions of a probability density, such as Shannon entropy and Rényi entropy. The density is assumed to be smooth, have bounded support, and be uniformly bounded from below on this set. Unlike previous $k$-NN estimators of non-linear density functionals, the proposed estimator uses data-splitting and boundary correction to achieve lower mean square error. Specifically, we assume that $T$ i.i.d. samples $\mathbf{X}_i \in \mathbb{R}^d$ from the density are split into two pieces of cardinality $M$ and $N$ respectively, with $M$ samples used for computing a k-nearest-neighbor density estimate and the remaining $N$ samples used for empirical estimation of the integral of the density functional. By studying the statistical properties of k-NN balls, explicit rates for the bias and variance of the BPI estimator are derived in terms of the sample size, the dimension of the samples and the underlying probability distribution. Based on these results, it is possible to specify optimal choice of tuning parameters $M/T$, $k$ for maximizing the rate of decrease of the mean square error (MSE). The resultant optimized BPI estimator converges faster and achieves lower mean squared error than previous $k$-NN entropy estimators. In addition, a central limit theorem is established for the BPI estimator that allows us to specify tight asymptotic confidence intervals.

## 1  Introduction

Non-linear functionals of a multivariate density $f$ of the form $\int g(f(x), x) f(x) dx$ arise in applications including machine learning, signal processing, mathematical statistics, and statistical communication theory. Important examples of such functionals include Shannon and Rényi entropy. Entropy based applications for image matching, image registration and texture classification are developed in [20, 34]. Entropy functional estimation is fundamental

to independent component analysis in signal processing [32]. Entropy has also been used in Internet anomaly detection [24] and data and image compression applications [23]. Several entropy based nonparametric statistical tests have been developed for testing statistical models including uniformity and normality [44, 10]. Parameter estimation methods based on entropy have been developed in [7, 37]. For further applications, see, for example, Leonenko *etal* [26].

In these applications, the functional of interest must be estimated empirically from sample realizations of the underlying densities. Several estimators of entropy measures have been proposed for general multivariate densities $f$. These include consistent estimators based on entropic graphs [19, 36], gap estimators [43], nearest neighbor distances [17, 26, 29, 45], kernel density plug-in estimators [1, 11, 3, 18, 4, 16], Edgeworth approximations [21], convex risk minimization [35] and orthogonal projections [25].

The class of density-plug-in estimators considered in this paper are based on $k$-nearest neighbor ($k$-NN) distances and, more specifically, bipartite k-nearest neighbor graphs over the random sample. The basic construction of the proposed bipartite plug-in (BPI) estimator is as follows (see Sec. II.A for a precise definition). Given a total of $T$ data samples we split the data into two parts of size $N$ and size $M$, $N + M = T$. On the part of size $M$ a $k$-NN density estimate is constructed. The density functional is then estimated by plugging the $k$-NN density estimate into the functional and approximating the integral by an empirical average over the remaining $N$ samples. This can be thought of as computing the estimator over a bipartite graph with the $M$ density estimation nodes connected to the $N$ integral approximating nodes. The BPI estimator exploits a close relation between density estimation and the geometry of proximity neighborhoods in the data sample. The BPI estimator is designed to automatically incorporate boundary correction, *without* requiring prior knowledge of the support of the density. Boundary correction compensates for bias due to distorted $k$-NN neighborhoods that occur for points near the boundary of the density support set. Furthermore, this boundary correction is *adaptive* in that we achieve the same MSE rate of convergence that can be attained using an oracle BPI estimator having knowledge of boundary of the support. Since the rate of convergence relates the number of samples $T = N+M$ to the performance of the estimator, convergence rates have great practical utility. A statistical analysis of the bias and variance, including rates of convergence, is presented for this class of boundary compensated BPI estimators. In addition, results on weak convergence (CLT) of BPI estimators are established. These results are applied to optimally select estimator tuning parameters $M/T, k$ and to derive confidence intervals. For arbitrary smooth functions $g$, we show that by choosing $k$ increasing in $T$ with order $O(T^{-2/(2+d)})$, an optimal MSE rate of order $O(T^{-4/(2+d)})$ is attained by the BPI estimator. For certain specific functions $g$ including Shannon entropy ($g(u) = \log(u)$) and Rényi entropy ($g(u) = u^{\alpha-1}$), a faster MSE rate of order $O(((\log T)^6/T)^{4/d})$ is achieved by BPI estimators by correcting for bias.

## 1.1   Previous work on $k$-NN functional estimation

The authors of [40, 17, 26, 29] propose $k$-NN estimators for Shannon entropy $(g(u) = \log(u))$ and Rényi entropy$(g(u) = u^{\alpha-1})$. Evans *etal* [13] consider positive moments of the $k$-NN distances $(g(u) = u^k, k \in \mathbb{N})$. Recently, Baryshnikov *etal* [2] proposed $k$-NN estimators for estimating $f$-divergence $\int \phi(f_0(x)/f(x))f(x)dx$ between an unknown density $f$, from which sample realizations are available, and a known density $f_0$. Because $f_0$ is known, the $f$-divergence $\int \phi(f_0(x)/f(x))f(x)dx$ is equivalent to a entropy functional $\int g(f(x), x)dx$ for a suitable choice of $g$. Wang *etal* [45] developed a $k$-NN based estimator of $\int g(f_1(x)/f_2(x), x)f_2(x)dx$ when both $f_1$ and $f_2$ are unknown. The authors of these works [40, 17, 13, 45] sestablish that the estimators they propose are asymptotically unbiased and consistent. The authors of [29] analyze estimator bias for $k$-NN estimation of Shannon and Rényi entropy. For smooth functions $g(.)$, Evans *etal* [12] show that the variance of the sums of these functionals of $k$-NN distances is bounded by the rate $O(k^5/T)$. Baryshnikov *etal* [2] improved on the results of Evans *etal* by determining the exact variance up to the leading term $(c_k/T$ for some constant $c_k$ which is a function of $k$). Furthermore, Baryshnikov *etal* show that the entropy estimator they propose converges weakly to a normal distribution. However, Baryshnikov *etal* do not analyze the bias of the estimators, nor do they show that the estimators they propose are consistent. Using the results obtained in this paper, we provide an expression for this bias in Section 4.4 and show that the optimal MSE for Baryshnikov's estimators is $O(T^{-2/(1+d)})$.

In contrast, the main contribution of this paper is the analysis of a general class of BPI estimators of smooth density functionals. We provide asymptotic bias and variance expressions and a central limit theorem. The bipartite nature of the BPI estimator enables us to correct for bias due to truncation of $k$-NN neighborhoods near the boundary of the support set; a correction that does not appear straightforward for previous $k$-NN based entropy estimators. We show that the BPI estimator is MSE consistent and that the MSE is guaranteed to converge to zero as $T \to \infty$ and $k \to \infty$ with a rate that is minimized for a specific choice of $k$, $M$ and $N$ as a function of $T$. Therefore, the thus optimized BPI estimator can be implemented without any tuning parameters. In addition a CLT is established that can be used to construct confidence intervals to empirically assess the quality of the BPI estimator. Finally, our method of proof is very general and it is likely that it can be extended to kernel density plug-in estimators, $f$-divergence estimation and mutual information estimation.

Another important distinction between the BPI estimator and the $k$-NN estimators of Shannon and Rényi entropy proposed by the authors of [40, 17, 26] is that these latter estimators are consistent for finite $k$, while the proposed BPI estimator requires the condition that $k \to \infty$ for MSE convergence. By allowing $k \to \infty$, the BPI estimators of Shannon and Rényi entropy achieve MSE rate of order $O(((\log T)^6/T)^{4/d})$. This asymptotic rate is faster than the $O(T^{-2/d})$ MSE convergence rate [29] of the previous $k$-NN estimators [40, 17, 26] that use a fixed value of $k$. It is shown by simulation that BPI's asymptotic performance advantages, predicted by our theory, also hold for small sample regimes.

## 1.2 Organization

The remainder of the paper is organized as follows. Section 2 formulates the entropy estimation problem and introduces the BPI estimator. The main results concerning the bias, variance and asymptotic distribution of these estimators are stated in Section 3 and the consequences of these results are discussed. The proofs are given in the Appendix. The MSE is analyzed in Section 4. We discuss bias correction of the BPI estimator for the case of Shannon and Rényi entropy estimation in Section 5. Estimation of Shannon MI is briefly discussed in Section 6. We numerically validate our theory by simulation in Section 7. Applications to structure discovery and dimension estimation are discussed in Sections 8 and 9 respectively. A conclusion is given in Section 10.

### Notation

Bold face type will indicate random variables and random vectors and regular type face will be used for non-random quantities. Denote the expectation operator by the symbol $\mathbb{E}$ and conditional expectation given $\mathbf{Z}$ by $\mathbb{E}_\mathbf{Z}$. Also define the variance operator as $\mathbb{V}[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2]$ and the covariance operator as $Cov[\mathbf{X}, \mathbf{Y}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])]$. Denote the bias of an estimator by $\mathbb{B}$.

## 2   Preliminaries

We are interested in estimating non-linear functionals $G(f)$ of $d$-dimensional multivariate densities $f$ with support $\mathcal{S}$, where $G(f)$ has the form

$$G(f) = \int g(f(x), x) f(x) d\mu(x) = \mathbb{E}[g(f(x), x)],$$

for some smooth function $g(f(x), x)$. Let $\mathcal{B}$ denote the boundary of $\mathcal{S}$. Here, $\mu$ denotes the Lebesgue measure and $\mathbb{E}$ denotes statistical expectation w.r.t density $f$. We assume that i.i.d realizations $\{\mathbf{X}_1, \ldots, \mathbf{X}_N, \mathbf{X}_{N+1}, \ldots, \mathbf{X}_{N+M}\}$ are available from the density $f$. Neither $f$ nor its support set are known.

The plug-in estimator is constructed using a data splitting approach as follows. The data is randomly subdivided into two parts $\mathcal{X}_N = \{\mathbf{X}_1, \ldots, \mathbf{X}_N\}$ and $\mathcal{X}_M = \{\mathbf{X}_{N+1}, \ldots, \mathbf{X}_{N+M}\}$ of $N$ and $M$ points respectively. In the first stage, a boundary compensated $k$-NN density estimator $\tilde{\mathbf{f}}_k$ is estimated at the $N$ points $\{\mathbf{X}_1, \ldots, \mathbf{X}_N\}$ using the $M$ realizations $\{\mathbf{X}_{N+1}, \ldots, \mathbf{X}_{N+M}\}$. Subsequently, the $N$ samples $\{\mathbf{X}_1, \ldots, \mathbf{X}_N\}$ are used to approximate the functional $G(f)$ to obtain the basic Bipartite Plug-In (BPI) estimator:

$$\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) = \frac{1}{N} \sum_{i=1}^{N} g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i). \tag{1}$$

As the above estimator performs an average over the $N$ variables $X_i$ of the function $g(\tilde{f}(X_i), X_i)$, which is estimated from the other $M$ variables, this estimator can be viewed as averaging over the edges of a bipartite graph with $N$ and $M$ nodes on its left and right parts.

## 2.1 Boundary compensated $k$-NN density estimator

Since the probability density $f$ is bounded above, the observations will lie strictly on the interior of the support set $\mathcal{S}$. However, some observations that occur close to the boundary of $\mathcal{S}$ will have $k$-NN balls that intersect the boundary. This leads to significant bias in the $k$-NN density estimator. In this section we describe a method that compensates for this bias. The method can be interpreted as extrapolating the location of the boundary from extreme points in the sample and suitably reducing the volumes of their $k$-NN balls.

Let $d(X, Y)$ denote the Euclidean distance between points $X$ and $Y$ and $\mathbf{d}_k(X)$ denote the Euclidean distance between a point X and its $k$-th nearest neighbor amongst the $M$ realizations $\mathbf{X}_{N+1}, .., \mathbf{X}_{N+M}$. Define a ball with radius $r$ centered at $X$: $S_r(X) = \{Y : d(X, Y) \leq r\}$. The $k$-NN region is $\mathbf{S}_k(X) = \{Y : d(X, Y) \leq \mathbf{d}_k(X)\}$ and the volume of the $k$-NN region is $\mathbf{V}_k(X) = \int_{\mathbf{S}_k(X)} dZ$. The standard $k$-NN density estimator [30] is defined as

$$\hat{\mathbf{f}}_k(X) = \frac{k-1}{M\mathbf{V}_k(X)}.$$

If a probability density function has bounded support, the $k$-NN balls $\mathbf{S}_k(X)$ centered at points $X$ close to the boundary may intersect with the boundary $\mathcal{B}$, or equivalently $\mathbf{S}_k(X) \cap \mathcal{S}^c \neq \phi$, where $\mathcal{S}^c$ is the complement of $\mathcal{S}$. As a consequence, the $k$-NN ball volume $\mathbf{V}_k(X)$ will tend to be higher for points $X$ close to the boundary leading to significant bias of the $k$-NN density estimator.

Let $R_k(X)$ correspond to the coverage value $(1+p_k)k/M$, i. e. , $R_k(X) = \inf\{r : \int_{S_r(X)} f(Z)dZ = (1 + p_k)k/M\}$, where $p_k = \sqrt{6}/(k^{\delta/2})$ for some fixed $\delta \in (2/3, 1)$. Define

$$\epsilon_{BC} = N \exp(-3k^{(1-\delta)}).$$

Define $N_k(X)$ as the region corresponding to the coverage value $(1 + p_k)k/M$, i.e. $N_k(X) = \{Y : d(X, Y) \leq R_k(X)\}$. Finally, define the interior region $\mathcal{S}_I$

$$\mathcal{S}_I = \{X \in \mathcal{S} : N_k(X) \cap \mathcal{S}^c = \phi\}. \tag{2}$$

We show in Appendix B that the bias of the standard $k$-NN density estimate is of order $O((k/M)^{(2/d)})$ for points $X \in \mathcal{S}_I$ and is of order $O(1)$ at points $X \in \mathcal{S} - \mathcal{S}_I$. This motivates the following method for compensating for this bias. This compensation is done in two stages: (i) the set of interior points $\mathcal{I}_N \subset \mathcal{X}_N$ are identified using variation in $k$-nearest neighbor distances in Algorithm 1 (see Appendix B for details) and it is show that $\mathcal{I}_N \notin \mathcal{S} - \mathcal{S}_I$ with probability $1 - O(\epsilon_{BC})$; and (ii) the density estimator at points in $\mathcal{B}_N = \mathcal{X}_N - \mathcal{I}_N$ are
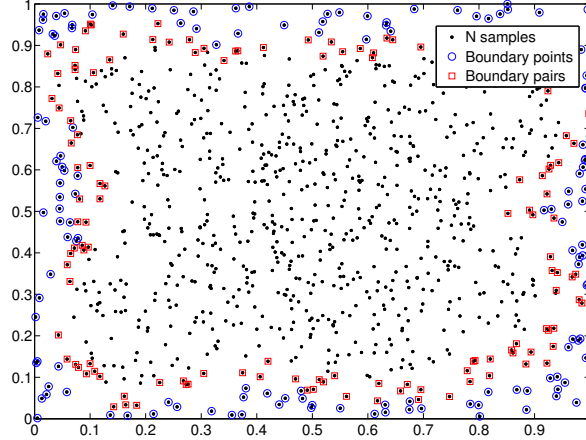
Figure 1: Detection of boundary points using Algorithm 1 for 2d beta distribution.

corrected by extrapolating to the density estimates at interior points $\mathcal{I}_N$ that are close to the boundary points. We emphasize that this nonparametric correction strategy does not assume knowledge about the support of the density $f$.

For each boundary point $\mathbf{X}_i \in \mathcal{B}_N$, let $\mathbf{X}_{n(i)} \in \mathcal{I}_N$ be the interior sample point that is closest to $\mathbf{X}_i$. The corrected density estimator $\tilde{\mathbf{f}}_k$ is defined as follows.

$$\tilde{\mathbf{f}}_k(\mathbf{X}_i) = \begin{cases} \hat{\mathbf{f}}_k(\mathbf{X}_i) & \{\mathbf{X}_i \in \mathcal{I}_N\} \\ \hat{\mathbf{f}}_k(\mathbf{X}_{n(i)}) & \{\mathbf{X}_i \in \mathcal{B}_N\} \end{cases} \tag{3}$$

# 3   Main results

Let $\mathbf{Z}$ denote an independent realization drawn from $f$. Also, define $\mathbf{Z}_{-1} \in \mathcal{S}_I$ to be $\mathbf{Z}_{-1} = \arg\min_{x \in \mathcal{S}_I} d(x, \mathbf{Z})$. Define $h(X) = \Gamma^{(2/d)}((d+2)/2)f^{-2/d}(X)tr[\nabla^2(f(X))]$. Denote the $n$-th partial derivative of $g(x,y)$ wrt $x$ by $g^{(n)}(x,y)$. Also, let $g'(x,y) := g^{(1)}(x,y)$ and $g''(x,y) := g^{(2)}(x,y)$. For some fixed $0 < \epsilon < 1$, define $p_l = ((k-1)/M)(1-\epsilon)\epsilon_0$ and $p_u = ((k-1)/M)(1+\epsilon)\epsilon_\infty$. Also define $\epsilon_1 = 1/(c_d \mathcal{D}^d)$, where $\mathcal{D}$ is the diameter of the bounded set $\mathcal{S}$ and define $q_l = ((k-1)/M)\epsilon_1$ and $q_u = (1+\epsilon)\epsilon_\infty$. Let $\mathbf{p}$ be a beta random variable with parameters $k, M-k+1$.

## 3.1   Assumptions

$(\mathcal{A}.0)$ : Assume that $M$, $N$ and $T$ are linearly related through the proportionality constant $\alpha_{frac}$ with: $0 < \alpha_{frac} < 1$, $M = \alpha_{frac}T$ and $N = (1 - \alpha_{frac})T$. $(\mathcal{A}.1)$ : Let the density $f$ be uniformly bounded away from 0 and finite on the set $\mathcal{S}$, i.e., there exist constants $\epsilon_0, \epsilon_\infty$ such
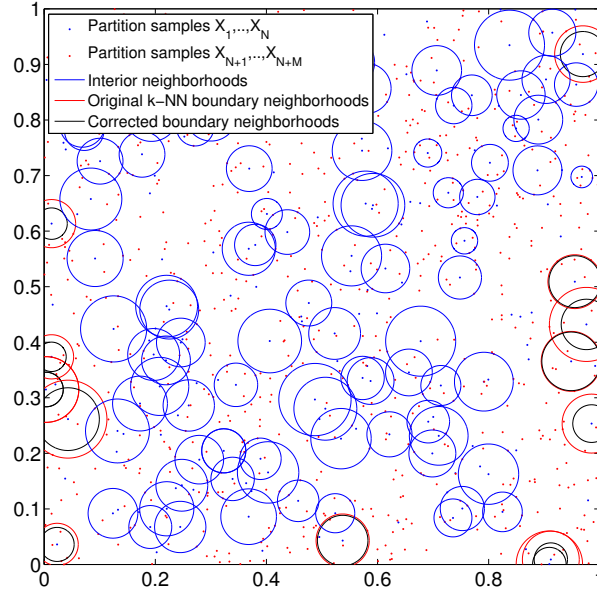
Figure 2: $k$-NN balls centered around a subsample of 2d uniformly distributed points. Note that the original $k$-NN balls centered at points close to boundary (red) over spill the boundary. The modified $k$-NN neighborhoods (black) corresponding to the corrected corrected density estimate $\tilde{\mathbf{f}}_k$ compensate for the over spill.

that $0 < \epsilon_0 \le f(x) \le \epsilon_\infty < \infty \; \forall x \in \mathcal{S}$. ($\mathcal{A}$.2): Assume that the density $f$ has continuous partial derivatives of order $2\nu$ in the interior of the set $\mathcal{S}$ where $\nu$ satisfies the condition $(k/M)^{2\nu/d} = o(1/M)$, and that these derivatives are upper bounded. ($\mathcal{A}$.3): Assume that the function $g(x, y)$ has $\lambda$ partial derivatives w.r.t. $x$, where $\lambda$ satisfies the conditions $k^{-\lambda} = o(1/M)$ and $O((\lambda^2((k/M)^{2/d} + 1/M))/M) = o(1/M)$. ($\mathcal{A}$.4): Assume that $\max\{6, 2\lambda\} < k <= M$. ($\mathcal{A}$.5): Assume that the absolute value of the functional $g(x, y)$ and its partial derivatives are strictly bounded away from $\infty$ in the range $\epsilon_0 < x < \epsilon_\infty$ for all $y$. ($\mathcal{A}$.6): Assume that $\sup_{x \in (q_l, q_u)} |(g^{(r)}/r!)^2(x, y)| e^{-3k^{(1-\delta)}} < \infty$, $\mathbb{E}[\sup_{x \in (p_l, p_u)} |(g^{(r)}/r!)^2(x/\mathbf{p}, y)|] < \infty$, for $r = 3, \lambda$.

## 3.2 Bias and Variance

Below the asymptotic bias and variance of the BPI estimator of general functionals of the density $f$ are specified. These asymptotic forms will be used to establish a form for the asymptotic MSE.

7

**Theorem 3.1.** *The bias of the BPI estimator* $\hat{\mathbf{G}}_k(f)$ *is given by*

$$\mathbb{B}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] \quad = \quad c_1 \left(\frac{k}{M}\right)^{2/d} + c_2 \left(\frac{1}{k}\right) + c_3(k, M, N) + O(\epsilon_{BC}) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right),$$

*where* $c_3(k, M, N) = \mathbb{E}[1_{\{\mathbf{Z} \in \mathcal{S} - \mathcal{S}_I\}}(g(f(\mathbf{Z}_{-1}), \mathbf{Z}_{-1}) - g(f(\mathbf{Z}), \mathbf{Z}))] = O(k/M)^{2/d}$, *and the constants* $c_1 = \mathbb{E}[g'(f(\mathbf{Z}), \mathbf{Z})h(\mathbf{Z})]$, $c_2 = \mathbb{E}[f^2(\mathbf{Z})g''(f(\mathbf{Z}), \mathbf{Z})/2]$.

**Theorem 3.2.** *The variance of the BPI estimator* $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ *is given by*

$$\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] \quad = \quad c_4 \left(\frac{1}{N}\right) + c_5 \left(\frac{1}{M}\right) + O(\epsilon_{BC}) + o\left(\frac{1}{M} + \frac{1}{N}\right),$$

*where the constants* $c_4 = \mathbb{V}[g(f(\mathbf{Z}), \mathbf{Z})]$ *and* $c_5 = \mathbb{V}[f(\mathbf{Z})g'(f(\mathbf{Z}), \mathbf{Z})]$.

*Proof.* We briefly sketch the proof here. The above theorems have been stated more generally and proved in Appendix D. The principal idea here involves Taylor series expansions of the functional $g(\tilde{\mathbf{f}}_k(X), X)$ about the true value $g(f(X), X)$, and subsequently (a) using the moment properties of density estimates derived in Appendix A to obtain the leading terms, and (b) bounding the remainder term in the Taylor series and showing that it can be ignored in comparison to the leading terms. $\qquad\square$

The leading terms $c_1(k/M)^{2/d} + c_2/k$ arise due to the bias and variance of $k$-NN density estimates respectively (see Appendix A), while the term $c_3(k, M, N)$ arises due to boundary correction (see Appendix B). Henceforth, we will refer to $c_3(k, M, N)$ by $c_3$. It is shown in Appendix B that $c_3 = O((k/M)^{2/d})$ (130). The term $O(\epsilon_{BC})$ arises from a concentration inequality that gives the probability of the event $\mathcal{I}_N \notin \mathcal{S} - \mathcal{S}_I$ as $1 - O(\epsilon_{BC})$. Observe that if $k$ increases logarithmically in $M$, specifically $(\log(M))^{2/(1-\delta)}/k \to 0$, then $O(\epsilon_{BC}) = o(N/M^3) = o(1/T)$.

The term $c_4/N$ is due to approximation of the integral $\int g(f(x), x)f(x)dx$ by the sample mean $(1/N)\sum_{i=1}^N g(f(\mathbf{X}_i), \mathbf{X}_i)$. The term $c_5/M$ on the other hand is due to the covariance between density estimates $\tilde{\mathbf{f}}(\mathbf{X}_i)$ and $\tilde{\mathbf{f}}(\mathbf{X}_j)$, $i \neq j$.

The constants $c_2, c_4$ and $c_5$ are once again functionals of the form $\int \tilde{g}(f(x), x)f(x)d\mu(x)$ and can be estimated using the proposed BPI estimator (1). On the other hand, the constant $c_1$ requires estimation of second order partial derivatives of $f$ in addition to estimating the density $f$. The partial derivatives might be estimated using the methods described in [38], $c_1$ could in principle be estimated in this manner.

To estimate $c_3$, we observe that $||\mathbf{Y} - \mathbf{Y}_{-1}|| = O((k/M)^{1/d})$ with probability $1 - O(N\mathcal{C}(k))$, and that $Pr(\mathbf{Y} \in \mathcal{S} - \mathcal{S}_I) = O((k/M)^{1/d})$. Let $h = \mathbf{Y} - \mathbf{Y}_{-1}$. Then,

$$\begin{aligned}
c_3 &= \mathbb{E}[1_{\{\mathbf{Y} \in \mathcal{S} - \mathcal{S}_I\}}(g(f(\mathbf{Y}_{-1}), \mathbf{Y}_{-1}) - g(f(\mathbf{Y}), \mathbf{Y}))] \\
&= \mathbb{E}[1_{\{\mathbf{Y} \in \mathcal{S} - \mathcal{S}_I\}}g'(f(\mathbf{Y}_{-1}), \mathbf{Y}_{-1})(f(\mathbf{Y}) - f(\mathbf{Y}_{-1}))] + O((k/M)^{3/d}) + O(\mathcal{C}(k)) \\
&= \mathbb{E}[1_{\{\mathbf{X}_1 \in \mathcal{S} - \mathcal{S}_I\}}g'(f(\mathbf{Y}_{-1}), \mathbf{Y}_{-1}) < \nabla f(\mathbf{Y}_{-1}), h >] + O((k/M)^{3/d}) + O(\mathcal{C}(k)).
\end{aligned}$$

The constant $c_3$ can then be estimated as

$$\hat{c}_3 = (1/N) \sum_{\mathbf{X}_i \in \mathcal{B}_N} g'(\hat{\mathbf{f}}_k(\mathbf{X}_{n(i)}), \mathbf{X}_{n(i)}) < \widehat{\nabla f}(\mathbf{X}_{n(i)}), \mathbf{X}_i - \mathbf{X}_{n(i)} >,$$

where the estimate $\widehat{\nabla} f$ of the gradient $\nabla f$ of $f$ might once again be estimated using the methods described in [38].

## 3.3 Central limit theorem

In addition to the results on bias and variance shown in the previous section, it is shown here that the BPI estimator, appropriately normalized, weakly converges to the normal distribution. The asymptotic behavior of the BPI estimator is studied under the following limiting conditions: (a) $k/M \to 0$, (b) $k \to \infty$ and (c) $N \to \infty$. As shorthand, the above limiting assumptions will be collectively denoted by $\Delta \to 0$.

**Theorem 3.3.** *The asymptotic distribution of the BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ is given by*

$$\lim_{\Delta \to 0} Pr\left( \frac{\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - \mathbb{E}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)]}{\sqrt{\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)]}} \le \alpha \right) = Pr(\mathbf{S} \le \alpha),$$

*where $\mathbf{S}$ is a standard normal random variable.*

*Proof.* Define the random variables $\{\mathbf{Y}_{M,i}; i = 1, \ldots, N\}$ for any fixed $M$

$$\mathbf{Y}_{M,i} = \frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)]}{\sqrt{\mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)]}},$$

The key idea here is to recognize that $\mathbf{Y}_{M,i}$ are exchangeable random variables. Blum et.al. [5] showed that for exchangeable 0 mean, unit variance random variables $\mathbf{Z_i}$, the sum $\mathbf{S}_N = \frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{Z}_i$ converges in distribution to $N(0,1)$ if and only if $Cov(\mathbf{Z}_1, \mathbf{Z}_2) = 0$ and $Cov(\mathbf{Z}_1^2, \mathbf{Z}_2^2) = 0$. In our case,

$$\begin{aligned} Cov(\mathbf{Y}_{M,i}, \mathbf{Y}_{M,j}) &= O(1/M), \\ Cov(\mathbf{Y}_{M,i}^2, \mathbf{Y}_{M,j}^2) &= O(1/M). \end{aligned}$$

As $M$ gets large, we then have that $Cov(\mathbf{Y}_{M,i}, \mathbf{Y}_{M,j}) \to 0$ and $Cov(\mathbf{Y}_{M,i}^2, \mathbf{Y}_{M,j}^2) \to 0$. We then extend the work by Blum et.al. to show that convergence in distribution to $N(0,1)$ holds in our case as both $N$ and $M$ get large. These ideas are rigorously treated in Appendix E. $\square$

The CLT for $k$-NN estimators of Rényi entropy was alluded to by Leonenko et.al. [17] by inferring from experimental results. Theorem 3.3 establishes the CLT for BPI estimators of arbitrary functionals, including Rényi entropy. This result allows one to define approximate finite sample confidence intervals on the estimated values of the functionals and define p-values .
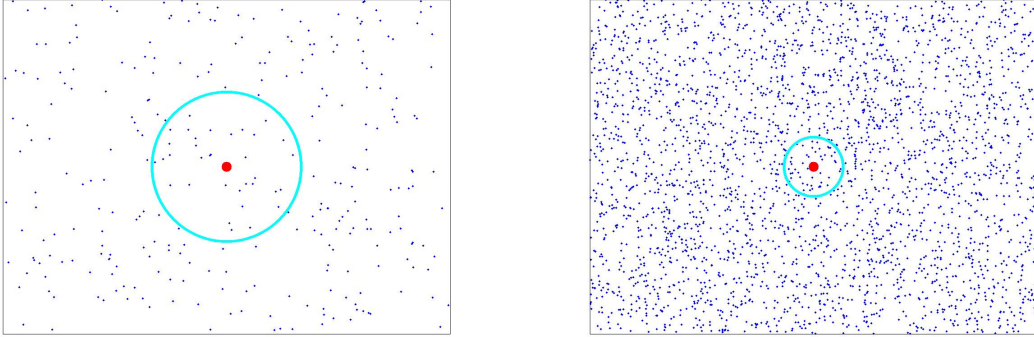
9

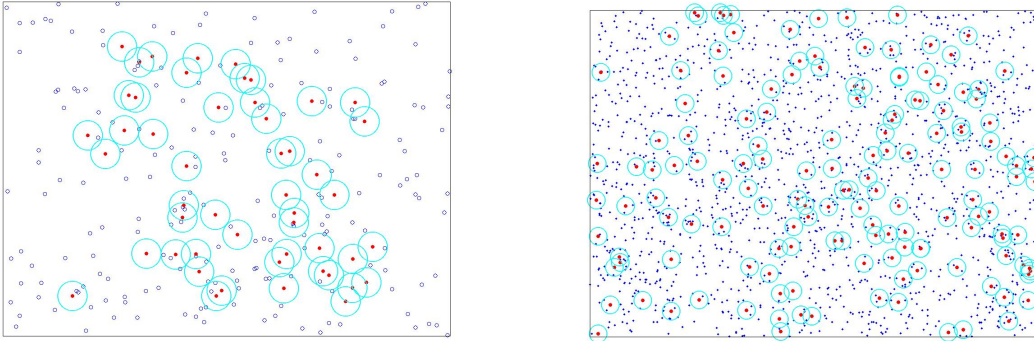Figure 3: Asymptotics. Variation of density estimate with increasing $k$ and $M$



Figure 4: Asymptotics. Variation of plug-in estimate with increasing $k$, $M$ and $N$

# 4 Analysis of M.S.E

Theorem 3.1 implies that $k \to \infty$ and $k/M \to 0$ in order that the BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ be asymptotically unbiased. Likewise, Theorem 3.2 implies that $N \to \infty$ and $M \to \infty$ in order that the variance of the estimator converge to 0. It is clear from Theorem 3.1 that the MSE is minimized when $k$ grows in polynomially in $M$. Throughout this section, we assume that $k = k_0 M^r$ for some $r \in (0,1)$. This implies that $O(\epsilon_{BC}) = O(N\mathcal{C}(k)) = o(1/M) = o(1/T)$. Figures 3 and 4 illustrate the asymptotic behavior of the density estimate and the plug-in estimate with increasing sample size.

## 4.1 Assumptions

Under the condition $k = k_0 M^r$, the assumptions ($\mathcal{A}$.2) and ($\mathcal{A}$.3) reduce to the following equivalent conditions: ($\mathcal{A}$.2): Let the density $f$ have continuous partial derivatives of order $2r$ in the interior of the set $\mathcal{S}$ where $r$ satisfies the condition $2r(1-t)/d > 1$. ($\mathcal{A}$.3): Let the functional $g(x, y)$ have $\lambda$ partial derivatives w.r.t. $x$, where $\lambda$ satisfies the conditions $t\lambda > 1$.

## 4.2 Optimal choice of parameters

In this section, we obtain optimal values for $k$, $M$ and $N$ for minimum M.S.E.

### 4.2.1 Optimal choice of $k$

Theorems III.1 and III.2 provide an optimal choice of $k$ that minimizes asymptotic MSE. Minimizing the MSE over $k$ is equivalent to minimizing the square of the bias over $k$. Define $c_o = c_1 + c_3/(k/M)^{2/d}$. The optimal choice of $k$ is given by

$$k_{opt} \quad = \quad \arg\min_k \mathbb{B}(\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)) = \lfloor k_0 M^{\frac{2}{2+d}} \rfloor, \tag{4}$$

where $\lfloor x \rfloor$ is the closest integer to $x$, and the constant $k_0$ is defined as $k_0 = (|c_2|d/2|c_0|)^{\frac{d}{d+2}}$ when $c_0 c_2 > 0$ and as $k_0 = (|c_2|/|c_0|)^{\frac{d}{d+2}}$ when $c_0 c_2 < 0$.

Observe that the constants $c_0$ and $c_2$ can possibly have opposite signs. When $c_0 c_2 > 0$, the bias evaluated at $k_{opt}$ is $b_0^+ M^{\frac{-2}{2+d}}(1+o(1))$ where $b_0^+ = c_0 k_0^{2/d} + c_2/k_0$. Let $k_{frac} = k_0 M^{\frac{2}{2+d}} - k_{opt}$. When $c_0 c_2 < 0$, observe that $c_0((k_{frac} + k_{opt})/M)^{2/d} + c_2/(k_{frac} + k_{opt})$ is equal to zero. When $c_0 c_2 < 0$, a higher order asymptotic analysis is required to specify the bias at the optimal value of $k$. In particular,

$$
\begin{aligned}
\mathbb{B}(\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)) \quad = \quad & c_1 \left( \frac{k}{M} \right)^{2/d} + c_2 \left( \frac{1}{k} \right) \\
& + h_1 \left( \frac{k}{M} \right)^{4/d} + h_2 \left( \frac{1}{k^2} \right) + h_3 \left( \left( \frac{k}{M} \right)^{2/d} \frac{1}{k} \right) \\
& + o \left( \left( \frac{k}{M} \right)^{4/d} + \frac{1}{k^2} + \left( \frac{k}{M} \right)^{2/d} \frac{1}{k} \right)
\end{aligned}
$$

where the constants are given by

$$h_1 = \mathbb{E}[(1/2)g''(f(\mathbf{Y}))h^2(X) + g'(f(\mathbf{Y}))h_o(\mathbf{Y})],$$

$$h_2 = \mathbb{E}[(2/3)g'''(f(\mathbf{Y}))f^3(\mathbf{Y})]$$

and

$$h_3 = (1 - 2/d)\mathbb{E}[g''(f(\mathbf{Y}))f(\mathbf{Y})c(\mathbf{Y})].$$

The bias evaluated at $k_{opt}$ is then given by $b_0^- M^{\frac{-4}{2+d}}(1 + o(1))$ where the constant $b_0^- = h_1 k_0^{4/d} + (h_2 + c_2 k_{frac})/k_0^2 + (h_3 + 2c_1 k_{frac}/d)k_0^{2/d-1}$.

Even though the optimal choice $k_{opt}$ depends on the unknown density $f$ (via the constant $k_0$), we observe from simulations that simply matching the rates, i.e. choosing $k = \bar{k} = M^{2/(2+d)}$, leads to significant MSE improvement. This is illustrated in Section 7.

### 4.2.2 Choice of $\alpha_{frac} = M/T$

Observe that the MSE of $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ is dominated by the squared bias $(O(M^{-4/(2+d)}))$ as contrasted to the variance $(O(1/N + 1/M))$. This implies that the MSE rate of convergence is invariant to the choice of $\alpha_{frac}$. This is corroborated by the experimental results shown in Fig. 12.

### 4.2.3 Discussion on optimal choice of $k$

The optimal choice of $k$ grows at a smaller rate as compared to the total number of samples $M$ used for the density estimation step. Furthermore, the rate at which $k/M$ grows decreases as the dimension $d$ increases. This can be explained by observing that the choice of $k$ primarily controls the bias of the entropy estimator. For a fixed choice of $k$ and $M$ ($k < M$), one expects the bias in the density estimates (and correspondingly in the estimates of the functional $G(f)$) to increase as the dimension increases. For increasing dimension an increasing number of the $M$ points will be near the boundary of the support set. This in turn requires choosing a smaller $k$ relative to $M$ as the dimension $d$ grows.

## 4.3 Optimal rate of convergence

Observe that the optimal bias decays as $b_0^+(T^{\frac{-2}{2+d}})(1 + o(1))$ when $c_0 c_2 > 0$ and $b_o^-(T^{\frac{-4}{2+d}})(1 + o(1))$ when $c_0 c_2 < 0$. The variance decays as $\Theta(1/T)(1 + o(1))$.

## 4.4 Comparison with results by Baryshnikov *etal*

Recently, Baryshnikov *etal* [2] have developed asymptotic convergence results for estimators of $f$-divergence $G(f_0, f) = \int f(x)\phi(f_0(x)/f(x))dx$ for the case where $f_0$ is known. Their estimators are based on sums of functionals of $k$-NN distances. They assume that they have $T$ i.i.d realizations from the unknown density $f$, and that $f$ and $f_0$ are bounded away from $0$ and $\infty$ on their support. The general form of the estimator of Baryshnikov *etal* is given by

$$\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS}) = \frac{1}{T}\sum_{i=1}^{T} g(\hat{\mathbf{f}}_{kS}(\mathbf{X}_i)),$$

where $\hat{\mathbf{f}}_{kS}(\mathbf{X}_i)$ is the standard $k$-NN density estimator [31] estimated using the $T-1$ samples $\{\mathbf{X}_1, .., \mathbf{X}_T\} - \{\mathbf{X}_i\}$.

Baryshnikov *etal* do not show that their estimator is consistent and do not analyze the bias of their estimator. They show that the leading term in the variance is given by $c_k/T$ for some constant $c_k$ which is a function of the number of nearest neighbors $k$. Finally they show that their estimator, when suitably normalized, is asymptotically normal. In contrast, we assume higher order conditions on continuity of the density $f$ and the functional $g$ (see Section 3) as compared to Baryshnikov *etal* and provide results on bias, variance and asymptotic distribution of data-split $k$-NN functional estimators of entropies of the form $G(f) = \int g(f(x))f(x)dx$. Note that we also require the assumption that $f$ is bounded away from 0 and $\infty$ on its support. Because we are able to establish expressions on both the bias and variance of the BPI estimator, we are able to specify optimal choice of free parameters $k, N, M$ for minimum MSE.

For estimating the functional $G(f) = \int g(f(x))f(x)dx$, the estimator of Baryshnikov can be used by restricting $f_0$ to be uniform. In Appendix C it is shown that under the additional assumption that $(\mathcal{A}.6)$ is satisfied by $\tilde{g} = g$, the bias of $\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})$ is

$$\mathbb{B}(\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})) = O((k/T)^{1/d}) + O(1/k). \tag{5}$$

In contrast, Theorem III. 1 establishes that the bias of the BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ decays as $\Theta((k/M)^{2/d} + 1/k) + O(\epsilon_{BC})$ and the variance decays as $\Theta(1/T)$. The bias of the BPI estimator has a higher exponent ($2/d$ as opposed to $1/d$) and this is a direct consequence of using the boundary compensated density estimator $\tilde{\mathbf{f}}_k$ in place of $\hat{\mathbf{f}}_k$.

It is clear from 5 that the estimator of Baryshnikov will be unbiased iff $k \to \infty$ as $T \to \infty$. Furthermore, the optimal rate of growth of $k$ is given by $k = T^{1/(1+d)}$. Furthermore, $c_k = \Theta(1)$ and therefore the overall optimal bias and variance of $\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})$ is given by $\Theta(T^{-1/(1+d)})$ and $\Theta(T^{-1})$ respectively. On the other hand, the optimal bias of the BPI estimator decays as $b_0^+(T^{\frac{-2}{2+d}})(1 + o(1))$ when $c_1c_2 > 0$ and $b_o^-(T^{\frac{-4}{2+d}})(1 + o(1))$ when $c_1c_2 < 0$ and the optimal variance decays as $\Theta(1/T)$. The BPI estimator therefore has faster rate of MSE convergence. Experimental MSE comparison of Baryshnikov's estimator against the proposed BPI estimator is shown in Fig. 12.

# 5 Bias correction factors

When the density functional of interest is the Shannon entropy ($g(u) = -\log(u)$) or the Rényi -$\alpha$ entropy($g(u) = u^{\alpha-1}$), a bias correction can be added to the BPI estimator that accelerates rate of convergence. Goria et.al. [26] and Leonenko et.al. [17] developed consistent Shannon and Rényi estimators with bias correction. The authors of [29] analyzed the bias for these estimators. When combined with the results of Baryshnikov *etal*, one can easily deduce the variance of these estimators and establish a CLT.

Let $\hat{\mathbf{H}}_S$ be the Shannon entropy estimate $\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})$ with the choice of functional $g(x) =$

$-\log(x)$. Let $\hat{\mathbf{I}}_{\alpha,S}$ be the estimate of the Rényi $\alpha$-integral estimate $\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})$ with the choice of functional $g(x) = x^{\alpha-1}$. Define $\tilde{\mathbf{H}}_S = \hat{\mathbf{H}}_S + [\log(k-1) - \Psi(k)]$, where $\psi(.)$ is the digamma function, and $\tilde{\mathbf{I}}_{\alpha,S} = [(\Gamma(k+(1-\alpha))/\Gamma(k))(k-1)^{\alpha-1}]^{-1}\hat{\mathbf{I}}_{\alpha,S}$. Also define the Rényi entropy estimator to be $\tilde{\mathbf{H}}_{\alpha,S} = (1-\alpha)^{-1}\log(\tilde{\mathbf{I}}_{\alpha,S})$. The estimators $\tilde{\mathbf{H}}_S$ and $\tilde{\mathbf{H}}_{\alpha,S}$ are the Shannon and Rényi entropy estimators of Goria *etal* [17] and Leonenko *etal* [26] respectively. In [29], it is shown that the bias of $\tilde{\mathbf{H}}_S$ and $\tilde{\mathbf{I}}_{\alpha,S}$ is given by $\Theta((k/T)^{1/d})$, while the variance was shown by Baryshnikov *etal* to be $O(1/T)$. In contrast, by (5), the bias of $\hat{\mathbf{H}}_S$ and $\hat{\mathbf{I}}_{\alpha,S}$ is given by $\Theta((k/T)^{1/d} + (1/k))$ (5). This can be understood as follows. From the results by [29], we have

$$\mathbb{E}[\hat{\mathbf{H}}_S] = I - [\log(k-1) - \Psi(k)] + c_{0,0}(k/T)^{1/d} + o((k/T)^{1/d}) \tag{6}$$

and

$$\mathbb{E}[\hat{\mathbf{I}}_{\alpha,S}] = [(\Gamma(k+(1-\alpha))/\Gamma(k))(k-1)^{\alpha-1}]I_\alpha + c_{0,\alpha}(k/T)^{1/d} + o((k/T)^{1/d}) \tag{7}$$

for some functionals of the density $c_{0,0}$ and $c_{0,\alpha}$. Note that $[(\Gamma(k+(1-\alpha))/\Gamma(k))(k-1)^{\alpha-1}] = 1 + O(1/k)$ and $\Psi(k) = \log(k-1) + O(1/k)$ as $k \to \infty$. From the above equations, the scale factor $[(\Gamma(k+(1-\alpha))/\Gamma(k))(k-1)^{\alpha-1}]$ and the additive factor $[\log(k-1) - \Psi(k)]$ account for the $O(1/k)$ terms in the expressions for bias of $\hat{\mathbf{H}}_S$ and $\hat{\mathbf{I}}_{\alpha,S}$, thereby removing the requirement that $k \to \infty$ for asymptotic unbiasedness. These bias corrections can be incorporated into the BPI estimator as follows.

## 5.1 Main results

For a general function $g(x,y)$, if there exist functions $g_1(k,M)$ and $g_2(k,M)$, such that

$$\begin{aligned}
&(i) \quad \mathbb{E}[g((k-1)x/M\mathbf{p},y)] = g(x,y)g_1(k,M) + g_2(k,M) + o(1/M), \\
&(ii) \quad ((k-1)/M)\mathbb{E}[g'((k-1)x/M\mathbf{p},y)\mathbf{p}^{2/d-1}] = g'(x,y)(k/M)^{2/d} + o((k/M)^{2/d}), \\
&(iii) \quad \lim_{k\to\infty} g_1(k,M) = 1, \\
&(iv) \quad \lim_{k\to\infty} g_2(k,M) = 0,
\end{aligned} \tag{8}$$

then define the BPI estimator with bias correction as

$$\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) \quad = \quad \frac{\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - g_2(k,M)}{g_1(k,M)}. \tag{9}$$

### 5.1.1 Bias and Variance

In addition to the assumptions listed in section 3.1, assume that $k = O((\log(M))^{2/(1-\delta)})$. Below the asymptotic bias and variance of the BPI estimator with bias correction are specified.

**Theorem 5.1.** *The bias of the BPI estimator* $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ *is given by*

$$\mathbb{B}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)] \quad = \quad c_1 \left(\frac{k}{M}\right)^{2/d} + c_3(k,M,N) + o\left(\left(\frac{k}{M}\right)^{2/d}\right). \tag{10}$$

14

**Theorem 5.2.** *The variance of the BPI estimator* $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ *is given by*

$$\mathbb{V}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)] \;\; = \;\; c_4\left(\frac{1}{N}\right) + c_5\left(\frac{1}{M}\right) + o\left(\frac{1}{M} + \frac{1}{N}\right).$$

### 5.1.2   CLT

**Theorem 5.3.** *The asymptotic distribution of the BPI estimator* $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ *is given by*

$$\lim_{\Delta \to 0} Pr\left(\frac{\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) - \mathbb{E}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)]}{\sqrt{\mathbb{V}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)]}} \le \alpha\right) = Pr(\mathbf{S} \le \alpha),$$

*where* $\mathbf{S}$ *is a standard normal random variable.*

### 5.1.3   MSE

Theorem IV. 1 specifies the bias of the BPI estimator, $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$, as $\Theta((k/M)^{2/d})$. Theorem IV. 2 specifies the variance as $\Theta(1/N + 1/M)$. By making $k$ increase logarithmically in $M$, specifically, $k = O((\log(M))^{2/(1-\delta)})$ for any value $\delta \in (2/3, 1)$, the MSE is given by the rate $\Theta(((\log(T))^{2/(1-\delta)}/T)^{4/d})$. The BPI estimator therefore has a faster rate of convergence in comparison to both Baryshnikov *etal*'s estimators $\hat{\mathbf{H}}_S$ and $\hat{\mathbf{I}}_{\alpha,S}$ (MSE $= \Theta(T^{-2/(1+d)})$) and Leonenko *etal*'s and Goria *etal*'s estimators $\tilde{\mathbf{H}}_S$ and $\tilde{\mathbf{I}}_{\alpha,S}$ (MSE $= \Theta(T^{-2/d})$). Experimental MSE comparison of Leonenko's estimator against the BPI estimator in Section V shows the MSE of the BPI estimator to be significantly lower. Finally, note that such bias correction cannot be applied for general entropy functionals, and the bias correction factors cannot in general be incorporated. In the next section, the application of BPI estimators for estimation of Shannon and Rényi entropies is illustrated.

## 5.2   Shannon and Rényi entropy estimation

For the case of Shannon entropy ($g(u) = -\log(u)$), it can be verified that $g_1(k, M) = 1$, $g_2(k, M) = \psi(k) - \log(k-1)$ satisfy (8). Similarly, for the case of Rényi entropy ($g(u) = u^{\alpha-1}$), $g_1(k, M) = (\Gamma(k)/\Gamma(k+1-\alpha))(1/(k-1)^{\alpha-1})$, $g_2(k, M) = 0$ satisfy (8).

For Shannon entropy ($g(u) = -\log(u)$) and Rényi entropy ($g(u) = u^{\alpha-1}$), the assumptions in Section 3.1 reduce to the following under the condition $k = O((\log(M))^{2/(1-\delta)})$. Assumption ($\mathcal{A}.1$) is unchanged. Assumption ($\mathcal{A}.2$) holds for any $r$ such that $2r > d$. The assumption ($\mathcal{A}.3$) is satisfied by the choice of $\lambda = \log(M)$. Assumption ($\mathcal{A}.4$) holds for ($g(u) = -\log(u)$) and ($g(u) = u^{\alpha-1}$). Next, it will be shown that ($\mathcal{A}.5$) is also satisfied by ($g(u) = -\log(u)$) and ($g(u) = u^{\alpha-1}$).

We note that $\tilde{g} = (g^{(3)}/6)^2$ for the choice of $g(u) = -\log(u)$ is given by $\tilde{g} = cu^{-6}$ for some constant $c$. Therefore,

$$
\begin{aligned}
\sup_{x \in (q_l, q_u)} |\tilde{g}(x,y)|e^{-3k^{(1-\delta)}} &= |c\epsilon_1^{-6}|(M/k)^6 O(e^{-3k^{(1-\delta)}}) \\
&= |c\epsilon_1^{-6}|(M/k)^6 O(e^{-3(\log(M))^2}) \\
&= |c\epsilon_1^{-6}|O(e^{-3(\log(M))^2 + 6\log(M) - 6\log(k)}) = o(1),
\end{aligned}
$$

and by (64), $\mathbb{E}[\sup_{x \in (p_l, p_u)} |\tilde{g}(x/\mathbf{p}, y)|] = |c|((1-\epsilon)\epsilon_0)^{-6}\mathbb{E}[(M\mathbf{p}/(k-1))^6] = |c|((1-\epsilon)\epsilon_0)^{-6}O(1) = O(1)$. Similarly, $\tilde{g} = (g^{(\lambda)}/(\lambda!))^2$ for the choice of $g(u) = -\log(u)$ is given by $\tilde{g} = \lambda^{-2}u^{-2\lambda}$. Then,

$$
\begin{aligned}
\sup_{x \in (q_l, q_u)} |\tilde{g}(x,y)|e^{-3k^{(1-\delta)}} &= O((M/k)^{2\lambda}e^{-3k^{(1-\delta)}}) \\
&= O((M/k)^{2\lambda}e^{-3(\log(M))^2}) \\
&= O(e^{-3(\log(M))^2 + 2(\log(M))^2 - 2\log(M)\log(k)}) = o(1),
\end{aligned}
$$

and by (64), $\mathbb{E}[\sup_{x \in (p_l, p_u)} |\tilde{g}(x/\mathbf{p}, y)|] = O(\mathbb{E}[(M\mathbf{p}/(k-1))^{2\lambda})]) = O(1)$. In an identical manner, $(\mathcal{A}.5)$ is satisfied when $g(u) = u^{\alpha-1}$.

To summarize, for functions $g(u) = -\log(u)$ and $g(u) = u^{\alpha-1}$, Theorem 5.1, 5.2 and 5.3 hold under the following assumptions: (i) $(\mathcal{A}.0)$, (ii) $(\mathcal{A}.1)$, (iii) the density $f$ has bounded continuous partial derivatives of order greater than $d$ and (iv) $k = O((\log(M))^{2/(1-\delta)})$. Furthermore the proposed BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ can be used to estimate Shannon entropy $(g(u) = -\log(u))$ and Rényi entropy $(g(u) = u^{\alpha-1})$ at MSE rate of $\Theta(((\log(T))^{2/(1-\delta)}/T)^{4/d})$.

# 6 Estimation of Shannon Mutual information

The joint entropy of random vectors $\mathbf{X}$ and $\mathbf{Y}$ with joint density $f_{XY}$ is given by

$$
H(\mathbf{X}, \mathbf{Y}) = -\int f_{XY} \log(f_{XY}) d\mu, \tag{11}
$$

where $f_{XY}$ is the joint density of $\mathbf{X}$ and $\mathbf{Y}$. The Shannon MI between two random vectors $\mathbf{X}$ and $\mathbf{Y}$ is then given by

$$
I(\mathbf{X}; \mathbf{Y}) = H(\mathbf{X}) + H(\mathbf{Y}) - H(\mathbf{X}, \mathbf{Y}). \tag{12}
$$

We use the following BPI estimator to estimate Shannon MI from $N + M$ $d$-dimensional i.i.d samples $\{(\mathbf{X_i}, \mathbf{Y_i}); i = 1, \ldots, N + M\}$ of the underlying joint density $f_{XY}$. We estimate the Shannon MI by estimating the individual entropies. We estimate the joint Shannon entropy $H(\mathbf{X}, \mathbf{Y})$ from samples using the *plug-in* estimate

$$
\hat{\mathbf{H}}(\mathbf{X}, \mathbf{Y}) = \frac{1}{N} \sum_{i=1}^{N} -\log(\tilde{\mathbf{f}}_\mathbf{k}(\mathbf{X_i}, \mathbf{Y_i})) + \log(k-1) - \psi(k), \tag{13}
$$

16

where $\hat{\mathbf{f}}_{\mathbf{XY}}$ is a $k$ nearest neighbor density estimate ($k$NN) estimated using the remaining $M$ samples.

The $k$NN density estimate [30] is given by

$$\tilde{\mathbf{f}}_{\mathbf{k}}(X, Y) = \frac{k-1}{M\mathbf{V}_{\mathbf{k}}(X, Y)}, \tag{14}$$

where $\mathbf{V}_{\mathbf{k}}(X, Y)$ is the volume corresponding to the $k$th nearest neighbor distance between the point of density estimation $(X, Y)$ and the $M$ i.i.d samples $\{(\mathbf{X_i}, \mathbf{Y_i}); i = N+1, \ldots, N+M\}$.

We estimate the marginal entropies by first obtaining estimates of the marginal density using $k$NN density estimates

$$\tilde{\mathbf{f}}_{\mathbf{k}}(X) = \frac{k-1}{M\mathbf{V}_{\mathbf{k}}(X)}, \tag{15}$$

where $\mathbf{V}_{\mathbf{k}}(X)$ is the volume corresponding to the $k$th nearest neighbor distance between the point of density estimation $X$ and the $M$ i.i.d samples $\{\mathbf{X_i}; i = N+1, \ldots, N+M\}$, and then plugging the estimated marginals into Eq. 16.

$$\hat{\mathbf{H}}(\mathbf{X}) = \frac{1}{N} \sum_{i=1}^{N} -\log(\tilde{\mathbf{f}}_{\mathbf{k}}(\mathbf{X_i})) + \log(k-1) - \psi(k). \tag{16}$$

Define the BPI estimator of Shannon MI:

$$\hat{\mathbf{I}}_N = \hat{\mathbf{H}}(\mathbf{X}) + \hat{\mathbf{H}}(\mathbf{Y}) - \hat{\mathbf{H}}(\mathbf{X}, \mathbf{Y}). \tag{17}$$

We make the following assumptions: (i) ($\mathcal{A}.0$), (ii) ($\mathcal{A}.1$), (iii) the density $f_{XY}$ has bounded continuous partial derivatives of order greater than $d$ and (iv) $k = O((\log(M))^{2/(1-\delta)})$. Note that the results here require cross moments between density estimates of the joint and marginal densities, which while not discussed in this report, can be obtained in exactly the same manner as computing cross moments between the same density.

**Theorem 6.1.** *The bias of the BPI estimator $\hat{\mathbf{I}}_N$ is given by*

$$\mathbb{B}[\hat{\mathbf{I}}_N] = c_1 \left(\frac{k}{M}\right)^{2/d} + c_3(k, M, N) + o\left(\left(\frac{k}{M}\right)^{2/d}\right). \tag{18}$$

**Theorem 6.2.** *The variance of the BPI estimator $\hat{\mathbf{I}}_N$ is given by*

$$\mathbb{V}[\hat{\mathbf{I}}_N] = c_4 \left(\frac{1}{N}\right) + c_5 \left(\frac{1}{M}\right) + o\left(\frac{1}{M} + \frac{1}{N}\right),$$

*where*

$$c_v = Var\left[\log\left(\frac{f_X(\mathbf{X})f_Y(\mathbf{Y})}{f_{XY}(\mathbf{X}, \mathbf{Y})}\right)\right].$$
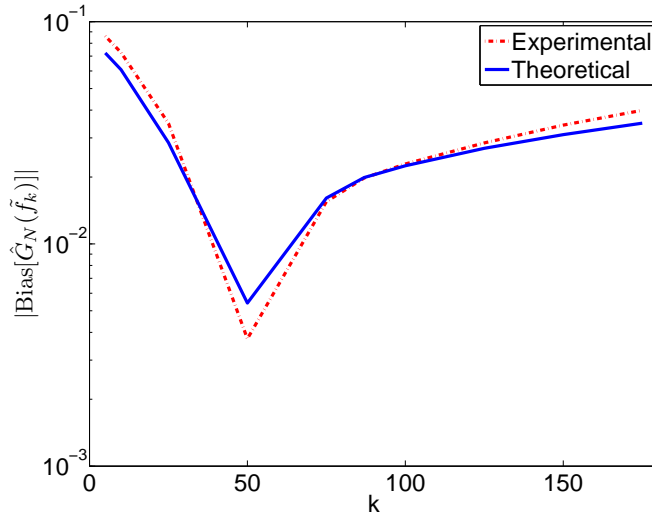
17

Figure 5: Comparison of theoretically predicted bias of BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ against experimentally observed bias as a function of $k$. The Shannon entropy $(g(u) = -\log(u))$ is estimated using the BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ on $T = 10^4$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). $N, M$ were fixed as $N = 3000$, $M = 7000$ respectively. The theoretically predicted bias agrees well with experimental observations. The predictions of our asymptotic theory therefore extend to the finite sample regime. The theoretically predicted optimal choice of $k_{opt} = 52$ also minimizes the empirical bias.

### 6.0.1  CLT

**Theorem 6.3.** *The asymptotic distribution of the BPI estimator $\hat{\mathbf{I}}_N$ is given by*

$$\lim_{\Delta \to 0} Pr\left(\frac{\hat{\mathbf{I}}_N - \mathbb{E}[\hat{\mathbf{I}}_N]}{\sqrt{\mathbb{V}[\hat{\mathbf{I}}_N]}} \leq \alpha\right) = Pr(\mathbf{S} \leq \alpha),$$

*where $\mathbf{S}$ is a standard normal random variable.*

# 7  Simulations

Here the theory established in Section 3 and Section 4 is validated. A three dimensional vector $\underline{X} = [X_1, X_2, X_3]^T$ was generated on the unit cube according to the i.i.d. Beta plus i.i.d. uniform mixture model:

$$f(x_1, x_2, x_3) = (1 - \epsilon) \prod_{i=1}^{3} f_{a,b}(x_i) + \epsilon, \tag{19}$$
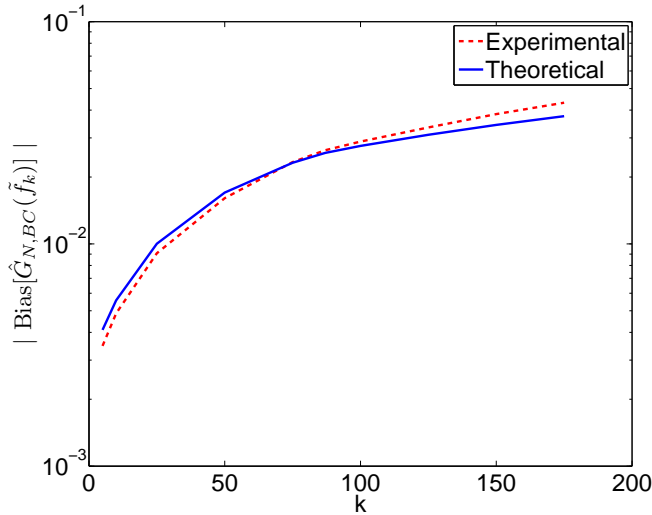
18

Figure 6: Comparison of theoretically predicted bias of BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ against experimentally observed bias as a function of $k$. The Shannon entropy ($g(u) = -\log(u)$) is estimated using the proposed BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ on $T = 10^4$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). $N, M$ were fixed as $N = 3000$, $M = 7000$ respectively. The empirical bias is in agreement with the bias approximations of Theorem IV. 1 and monotonically increases with $k$.

where $f_{a,b}(x)$ is a univariate Beta density with shape parameters $a$ and $b$. For the experiments the parameters were set to $a = 4, b = 4$, and $\epsilon = 0.2$. The Shannon entropy ($g(u) = -\log(u)$) is estimated using the BPI estimators $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ and $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$.

In Fig. 5, the bias approximations of Theorem III. 1 are compared to the empirically determined estimator bias of $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$. $N$ and $M$ are fixed as $N = 3000$, $M = 7000$. Note that the theoretically predicted optimal choice of $k_{opt} = 52$ minimizes the experimentally obtained bias curve. Thus, even though our theory is asymptotic it provides useful predictions for the case of finite sample size, specifying bandwidth parameters that achieve minimum bias. Further note that by matching rates, i.e. choosing $k = \bar{k} = M^{2/(2+d)} = 83$ also results in significantly lower MSE when compared to choosing $k$ arbitrarily ($k < 10$ or $k > 150$). In Fig. 6, the bias approximations of Theorem IV. 1 are compared to the empirically determined estimator bias of $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$. Observe that the empirical bias, in agreement with the bias approximations of Theorem IV. 1, monotonically increases with $k$.

In Fig. 7, the empirically determined variance of $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ is compared with the variance expressed by Theorem III. 2 for varying choices of $N$ and $M$, with fixed $N + M = 10,000$. The theoretically predicted variance agrees well with experimental observations. A Q-Q plot of the normalized BPI estimate $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ and the standard normal distribution is shown in Fig. 8. The linear Q-Q plot validates the Central Limit Theorem III. 3 on the uncompensated BPI estimator. To verify that the predicted confidence intervals were indeed as advertised, the empirically determined and theoretically predicted confidence intervals were compared
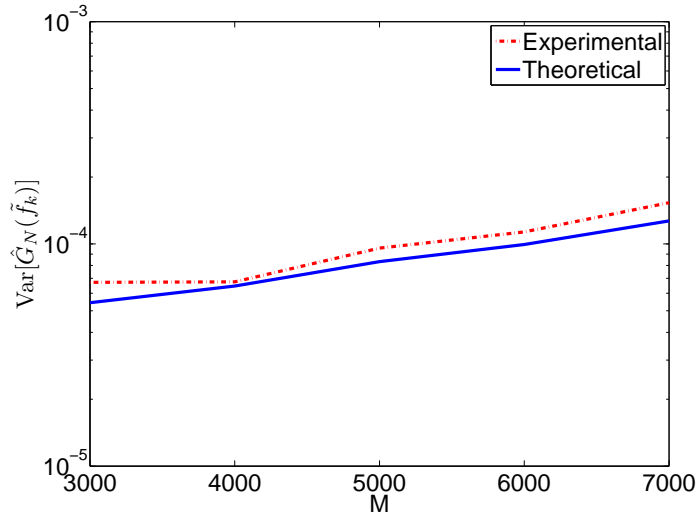
Figure 7: Comparison of theoretically predicted variance of BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ against experimentally observed variance as a function of $M$. The Shannon entropy $(g(u) = -\log(u))$ is estimated using the proposed BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ on $T = 10^4$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). $k$ is chosen to be $k_{opt} = k_0 M^{2/(2+d)}$. The theoretically predicted variance agrees well with experimental observations.

in Fig. 10. The lengths of the predicted confidence intervals are accurate to within 12% of the length of the true confidence intervals.

We additionally show in Fig. 11 a plot of the empirically determined estimator bias (via simulation) vs the bias predicted by our theory as a function of sample size $T$, which matches the theoretical prediction.

For Shannon entropy $(g(u) = -\log(u))$, the uncompensated and compensated BPI estimators are related by

$$\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) = \hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) + \log(k - 1) - \psi(k).$$

The variance and normalized distribution of these estimators are therefore identical. Consequently, Fig. 7 and Fig. 8 also validate Theorem IV. 2 and Theorem IV. 3 respectively.

Finally, using the CLT, the 95% coverage intervals of the BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ are shown as a function of sample size $T$ in Fig. 9. The lengths of the predicted confidence intervals are accurate to within 12% of the true confidence intervals (determined by simulation over the range of 80% to 100% coverage - data not shown). These coverage intervals can be interpreted as confidence intervals on the true entropy, provided that the constants $c_1, .., c_5$ can be accurately estimated.
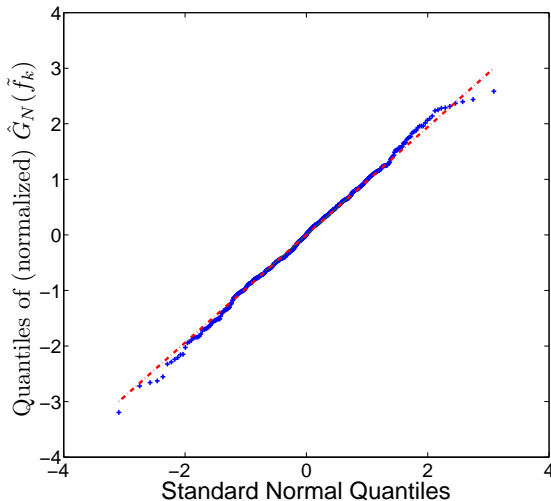
20

Figure 8: Q-Q plot comparing the quantiles of the BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ (with $g(u) = -\log(u)$) on the vertical axis to a standard normal population on the horizontal axis. The Shannon entropy $(g(u) = -\log(u))$ is estimated using the proposed BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ on $T = 10^4$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). $k, N, M$ are fixed as $k = k_{opt} = 52$, $N = 3000$ and $M = 7000$ respectively. The approximate linearity of the points validates our central limit theorem 3.3.

## 7.1 Experimental comparison of estimators

The Rényi $\alpha$-entropy $(g(u) = u^{\alpha-1})$ is estimated for $\alpha = 0.5$, with the same underlying 3 dimensional mixture of the beta and uniform densities defined above. Several estimators are compared: Baryshnikov's estimator $\hat{\mathbf{I}}_{\alpha,S}$, the $k$-NN estimator $\tilde{\mathbf{I}}_{\alpha,S}$ of Leonenko *etal* [17], the BPI estimator without bias correction $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ and the proposed BPI estimator with bias correction $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$. The results are shown in Fig. 12. It is clear from the figure that the BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ has the fastest rate of convergence, consistent with our theory. Note that, in agreement with our analysis in Section 4.4, the bias uncompensated BPI estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ outperforms Baryshnikov's estimator $\hat{\mathbf{I}}_{\alpha,S}$.

# 8 Application to structure discovery

Discovering structural dependencies among random variables from a multivariate sample is an important task in signal processing, pattern recognition and machine learning. Based on dependence relationships, the density function of the variables can be modeled using factor graphs. When the sample is highly structured, the corresponding factor graph configuration is sparse. Sparse factor graphs correspond to joint multivariate distributions which separate into a parsimonious product of few lower dimensional distributions. The inherent low-dimensional
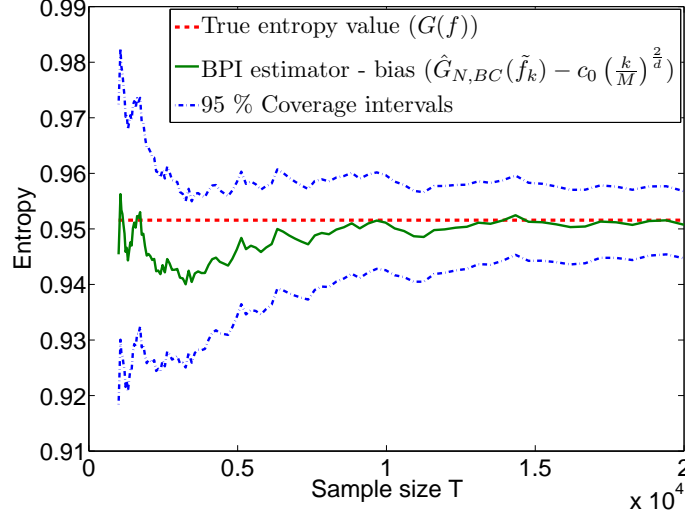
21

Figure 9: 95% coverage intervals of BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$, predicted using the Central limit theorem 3.3, as a function of sample size $T$. The Shannon entropy $(g(u) = -\log(u))$ is estimated using the proposed BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ on $T$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). The lengths of the coverage intervals are accurate to within 12% of the empirical confidence intervals obtained from the empirical distribution of the BPI estimator.
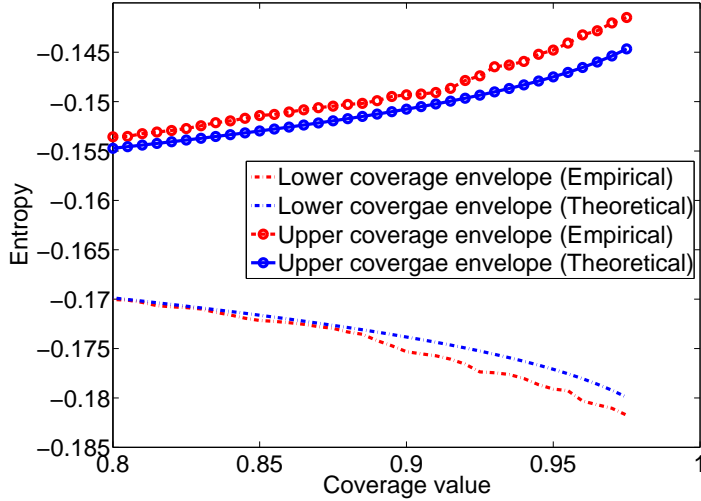


Figure 10: Empirically determined and theoretically predicted coverage envelopes as a function of coverage values. There is good agreement between the theoretically predicted and empirical coverage intervals.

nature of this product leads to a compact representation of the variables having sparse factor graph configurations.
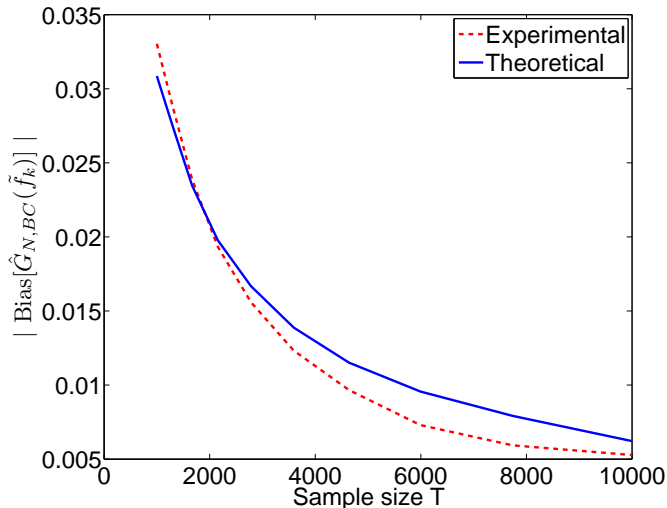
Figure 11: Comparison of theoretically predicted bias of BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ against experimentally observed bias as a function of sample size $T$. The Shannon entropy ($g(u) = -\log(u)$) is estimated using the proposed BPI estimator $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ on $T$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). The empirical bias is in agreement with the bias approximations of Theorem IV. 1 and monotonically decreases with $T$.

In practice, these structure dependencies have to be discovered from sample realizations of the multivariate distribution. Discovering dependencies when parametric probability density models are not known a priori is an important restriction of the above problem. For parametric distribution estimates, the errors are of order $O(1/N)$ if the true distribution is included in the parametric model. If not, a non-vanishing bias will dominate the error yielding an even higher error than that of a nonparametric distribution estimate (e.g. $k$NN estimates). In this restricted setting, recourse is therefore taken to nonparametric methods.

Chow et.al. [8] proposed an elegant solution to structure discovery of Markov tree distributions and provided a nonparametric algorithm to obtain the optimal tree. Ihler et.al. [22] developed the method of nonparametric hypothesis tests for structure discovery.

Nonparametric methods, while asymptotically consistent, can uncover incorrect factor graph structure when estimated from a finite number of samples. This is distinctly true for small sample sizes. While consistency is an important qualitative property, there is clearly an important motivation for quantitative characterization of performance in structure discovery. In this work, we analyze factor graph structure discovery in the finite sample size setting.

We present a class of $k$-nearest neighbor ($k$NN) based nonparametric geometric algorithms to discover factor graph structure among variables. We provide results on mean square error of the nonparametric estimates, which can be optimized over free parameters, thereby guaranteeing improved correct structure discovery. In addition, we provide confidence intervals
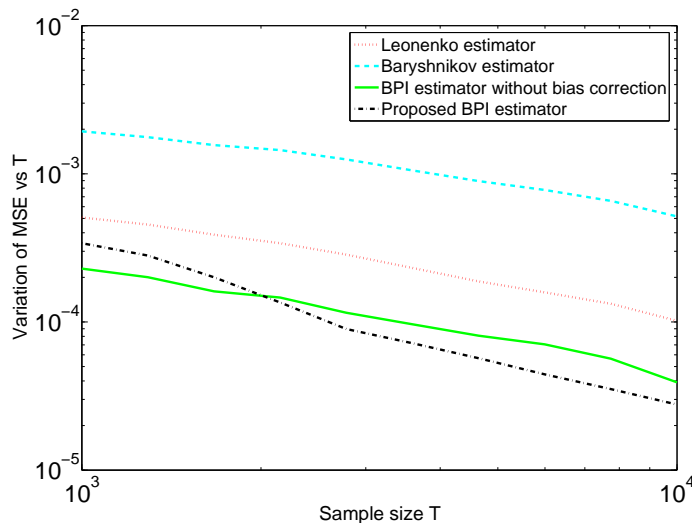
Figure 12: Variation of MSE of $k$-nearest neighbor estimator of Leonenko *etal* [17] and the $k$-nearest neighbor estimator of Baryshnikov *etal* [2] and BPI estimators with and without boundary correction, as a function of sample size $T$. The Rényi entropy $(g(u) = u^{\alpha-1})$ is estimated for $\alpha = 0.5$ using these estimators on $T$ i. i. d. samples drawn from the $d = 3$ dimensional uniform-beta mixture density (19). The figure shows that the proposed BPI estimator has the fastest rate of convergence.

on these nonparametric estimates to determine the probability of false error in choosing an incorrect structure model. These results are an direct extension of our work on optimized nonparametric estimates of divergence measures introduced earlier.

As a consequence of our statistical analysis, we introduce the notion of dependence-based **dimension** for factor graph models and show that comparing models within the same dimension class is an easier task with lower probability of false error as compared to comparing models across different dimensions.

## 8.1 Factor graphs

Factor graphs are bipartite graphs used to represent factorizations of probability density functions. Consider a set of variables $\underline{X} = \{X_1, X_2, \ldots, X_T\}$ and let $\{S_j \subseteq \{X_1, X_2, \ldots, X_n\}, j = 1, \ldots, m\}$ be a set of subsets of $\underline{X}$. Let $g(X_1, \ldots, X_T)$ denote a probability density function on the random vector $\underline{X}$. For the factorization $g(X_1, \ldots, X_T) = \prod_{j=1}^{m} f_j(S_j)$ of the density function, the corresponding factor graph $G = (\underline{X}, \underline{F}, E)$ consists of variable vertice's $\underline{X}$ , factor vertices's $\underline{F} = \{f_1, f_2, \ldots, f_m\}$, and edges $E$. The edges in the factor graph depend on the factorization as follows: there is an undirected edge between factor vertex $f_j$ and variable vertex $X_k$ when $X_k \subseteq S_j$.

## 8.2 Factor graph discovery

**Problem statement:** Consider a set of factor graphs $\{g_i(X_1, \ldots, X_T), i = 1, \ldots, I\}$. We seek to find the factor graph configuration from this set that best models the data.

The Kullback-Leibler (KL) divergence measure induces a **geometry** on the space of probability distributions. On this induced geometry, we naturally define the best factor graph configuration $g_o$ to be the one closest to the actual distribution $p(X_1, \ldots, X_T)$ in terms of KL divergence (c.f. [8]).

$$g_o = \arg\min_{g_i} KL(p||g_i) = \arg\min_{g_i} H_c(p, g_i), \tag{20}$$

where $H_c(p, g_i) = -\int p \log g_i$ is the cross-entropy between $p$ and $g_i$. In practice, these cross-entropy terms have to be estimated from the finite data sample. **Errors in estimation of cross-entropy terms can result in incorrect factor graph discovery**.

The problem considered by [8] is a specific instance of discovering factor graph structure. For the class of Markov tree factor graphs considered by [8], the cross entropy reduces to a sum of pairwise Shannon mutual information terms between variables with edges in the Markov tree. In their work, they empirically estimate the mutual information terms from the data using nonparametric estimators which are consistent. However, they do not take into account the error in the mutual information estimates when estimated from finite samples.

## 8.3 Disjoint factor graph discovery

In order to illustrate the effect of nonparametric estimation from finite sample size on factor graph discovery, we restrict our attention to disjoint factor graphs ([22]). For $i = 1, \ldots, I$, let

$$g_i(X_1, X_2, \ldots, X_T) = \prod_{j=1}^{m} p(S_j^{(i)}), \tag{21}$$

where $S_j^{(i)} \cap S_k^{(i)} = \phi$ whenever $j \neq k$, and $p(.)$ denotes the marginal density function. In this case of disjoint factor graphs, the cross-entropy takes the following simple form:

$$H_c(p, g_i) = \sum_j H(S_j^{(i)}), \tag{22}$$

where $H(S_j^{(i)})$ is the Shannon entropy of the variables $S_j^{(i)}$ under the true distribution $p$.

For example, consider the disjoint factor graph $g(X_1, \ldots, X_5) = p(X_1, X_2)p(X_3)p(X_4, X_5)$. The cross-entropy for this factor graph is given by $H_c(p, g) = H(X_1, X_2) + H(X_3) + H(X_4, X_5)$.

Consider two disjoint factor graph configurations: (a) $n(X_1, \ldots, X_T) = \prod_{i=1}^{m_1} f(R_i)$ and (b) $l(X_1, \ldots, X_T) = \prod_{j=1}^{m_2} f(S_j)$. Denote the dimension of $R_i$ by $d_i^n$ and $S_j$ by $d_j^l$. We note that $\sum_{i=1}^{m_1} d_i^{(n)} = \sum_{j=1}^{m_2} d_j^{(l)} = T$. Based on the above formulation, in order to compare the two potential factor graph models $n$ and $l$, we need to compare the respective cross-entropy terms. The cross entropy test is stated below.

**Cross entropy test:** The cross entropy test to compare between models $n$ and $l$ is given by

$$H_c(p,n) - H_c(p,l) = \sum_{i=1}^{m_1} H(R_i) - \sum_{j=1}^{m_2} H(S_j) \gtrless 0. \tag{23}$$

We estimate these entropy terms in the test statistic $H_c(p,n) - H_c(p,l)$ from sample realizations using $k$NN plug-in estimators introduced earlier.

## 8.4 Errors in factor graph discovery

To illustrate the effect of estimation error in factor graph discovery, again consider the two factor graph models $n(X_1, \ldots, X_T) = \prod_{i=1}^{m_1} f(R_i)$ and $l(X_1, \ldots, X_T) = \prod_{j=1}^{m_2} f(S_j)$.

The cross entropy test (Eq. 22) between models $n$ and $l$ is $H_c(p,n) - H_c(p,l) \gtrless 0$. We replace this optimal cross entropy test with the following **surrogate** cross entropy test:

$$\hat{H}_c(p,n) - \hat{H}_c(p,l) = \sum_{i=1}^{m_1} \hat{H}(R_i) - \sum_{j=1}^{m_2} \hat{H}(S_j) \gtrless 0. \tag{24}$$

where we estimate entropy terms $\hat{H}(R_i)$ or $\hat{H}(S_j)$ using independent realizations of the underlying density $p$. To elaborate, if we have $V$ samples $\{\underline{X}^{(1)}, \ldots, \underline{X}^{(V)}\}$ from the density $p$, we partition these $V$ samples into $m_1 + m_2$ disjoint subsets of size $N + M$ each. This implies that $N + M \approx V/(m_1 + m_2)$. We then use each subset to estimate entropy using the partitioning strategy as discussed earlier.

Denote the coefficients corresponding to the entropy estimate $\hat{H}(R_i)$ of the subset of variables $R_i$ in the factor graph model $n$ by $c_{n_i1}$, $c_{n_i2}$ and $c_{n_i4}$. Using the theorems established in this report, we have the following results:

**Mean:** The mean of this surrogate test statistic is then given by

$$
\begin{aligned}
\mathbb{E}_p[\hat{H}_c(p,n) - \hat{H}_c(p,l)] \;=\;\; & H_c(p,n) - H_c(p,l) \\
& + \sum_{i=1}^{m_1} c_{n_i1} \left(\frac{k}{M}\right)^{2/d_i^{(n)}} - \sum_{j=1}^{m_2} c_{l_j1} \left(\frac{k}{M}\right)^{2/d_j^{(l)}} \\
& + \sum_{i=1}^{m_1} c_{n_i2}/k - \sum_{j=1}^{m_2} c_{l_j2}/k. 
\end{aligned}
\tag{25}
$$

**Variance:** The variance of the surrogate test statistic is then given by the sum of the variance of the individual entropy estimates (by independence)

$$\mathbb{V}_p[\hat{H}_c(p,n) - \hat{H}_c(p,l)] \;=\; \left(\sum_{i=1}^{m_1} c_{n_i4} + \sum_{j=1}^{m_2} c_{l_j4}\right)\left(\frac{1}{N}\right). \tag{26}$$

**Weak convergence:** Again, by independence of the individual entropy estimates, we have the following weak convergence law

$$\lim_{N,M\to\infty} Pr\left(\frac{\sqrt{N}(\hat{H}_c(p,n) - \hat{H}_c(p,l) - \mathbb{E}_p[\hat{H}_c(p,n) - \hat{H}_c(p,l)])}{\sqrt{\mathbb{V}_p[\hat{H}_c(p,n) - \hat{H}_c(p,l)]}} \le \alpha\right) = Pr\left(Z \le \alpha\right), \quad (27)$$

where $Z$ is standard normal.

## 8.5 Discussion

From the above expressions for the mean, variance and weak convergence law of the surrogate test statistic, we make the following observations:

1. The bias term is dependent on the dimension of the factors of the factor graph models $d_i^{(n)}$ and $d_j^{(l)}$. The variance term is independent of dimension. Furthermore, it is clear that the bias term dominates the MSE as the dimension of the factors grows.

2. For better performance in discovering factor graph structure using cross entropy tests, it is clear that we want the MSE of the surrogate test statistic to be small. A significant route to achieving this is to get the bias from each factor graph cross entropy estimate in the estimated test statistic to cancel. This is to say, we want

$$\begin{aligned}
\mathbb{E}_p[\hat{H}_c(p,n) - \hat{H}_c(p,l)] &\approx H_c(p,n) - H_c(p,l) \\
\Rightarrow \mathbb{E}_p[\hat{H}_c(p,n)] - \hat{H}_c(p,n) &\approx \mathbb{E}_p[\hat{H}_c(p,l)] - \hat{H}_c(p,l) \\
\Rightarrow \sum_{i=1}^{m_1} c_{n_i 1}\left(\frac{k}{M}\right)^{2/d_i^{(n)}} + \sum_{i=1}^{m_1} c_{n_i 2}/k &\approx \sum_{j=1}^{m_2} c_{l_j 1}\left(\frac{k}{M}\right)^{2/d_j^{(l)}} + \sum_{j=1}^{m_2} c_{l_j 2}/k. \quad (28)
\end{aligned}$$

3. This cancellation effect will be maximized when the dimensions of the factor graph subsets $R_i$ and $S_j$ match. That is to say, we want $m_1 = m_2$ and furthermore $d_i^{(n)} = d_j^{(l)}$. In this case, the bias from each cross entropy estimate are of the same order and will nearly cancel.

   On the other hand, when there is a mismatch in dimension, the bias from one cross entropy estimate will dominate the bias from the other cross entropy estimate, resulting in significant bias in the surrogate test statistic.

   In both these cases, the variance of the surrogate test statistic will be of the same order $O(1/N)$.

4. This gives rise to notion of multivariate dimension for factor graphs. Index the factorizations according to the vector $E = [e_1, e_2, ..., e_p]$, where $e_i$ is an integer between $0$ and $T$ that counts the number of factors of order $i$, i.e. involving a marginal density

over $i$ variables. The **dimension** $E$ of factor graph configurations partitions the factor graphs into equivalence classes having nearly constant cross entropy estimate bias.

For two factor graph models $n$ and $l$ with dimensions $E_n$ and $E_l$, we will refer to $n$ as a higher dimensional model relative to $l$ if the last non-zero entry of $E_n - E_l$ is positive.

5. As discussed earlier, the bias will not be a significant factor when comparing models over an equivalence class having fixed values of $E$. On the other hand, the bias will be significant when comparing models across different values of $E$, resulting in higher probability of error in factor graph discovery.

6. Prior knowledge of the equivalence class will therefore translate into much improved performance in factor graph discovery as compared to prior knowledge that mixes between equivalence classes.

7. We note that the number of samples required to maintain a constant level of bias grows **geometrically** with dimension $E$.

8. Using the expressions for the bias and variance of the surrrogate test statistic, we can optimize over the free parameters: (a) the choice of partition $N$ and $M$ for fixed total sample size $N + M$ and (b) the choice of bandwidth parameter $k$, for minimum MSE.

9. Using the weak convergence law, we can theoretically predict the probability of choosing model $n$ over model $l$ using the surrogate cross entropy test.


## 8.6   Experiment

We illustrate the implications of our analysis with a toy example. Let $f_\beta(x, a, b, d)$ denote a beta density of dimension $d$ with parameters $a$ and $b$. Now let $f_\mu(x, d) = 0.5 f_\beta(x, 5, 2, d) + 0.5 f_\beta(x, 2, 5, d)$ be a mixture of beta densities. When $d > 1$, the mixing of densities ensures there is strong dependence between the variates.

We draw $V = 10^5$ independent sample realizations from the joint density $p(X_1, \ldots, X_5) = f_\mu(X_1, 1) f_\mu(X_2, 1) f_\mu(X_3, 1) f_\mu(X_4, X_5, 2)$.

|   | E | True | False |
|---|---|------|-------|
| l | $[1, 0, 0, 1, 0]$ | $f(X_1, X_2, X_4, X_5) f(X_3)$ | $f(X_1, X_2, X_3, X_4) f(X_5)$ |
| m | $[1, 2, 0, 0, 0]$ | $f(X_1, X_2) f(X_4, X_5) f(X_3)$ | $f(X_1, X_3) f(X_2, X_4) f(X_5)$ |
| n | $[3, 1, 0, 0, 0]$ | $f(X_4, X_5) f(X_1) f(X_2) f(X_3)$ | $f(X_2, X_4) f(X_1) f(X_3) f(X_5)$ |

**Experiment** The table above shows six different factor graph models. We compare each true model against each false model. Denote the true models by $l_T$, $m_T$ and $n_T$ and the corresponding false models by $l_F$, $m_F$ and $n_F$. We note that the true cross entropy terms $H_c(p, l_T) = H_c(p, m_T) = H_c(p, n_T)$ and $H_c(p, l_L) = H_c(p, m_L) = H_c(p, n_L)$. This guarantees level playing field when comparing each true model against each false model using the surrogate cross entropy test.
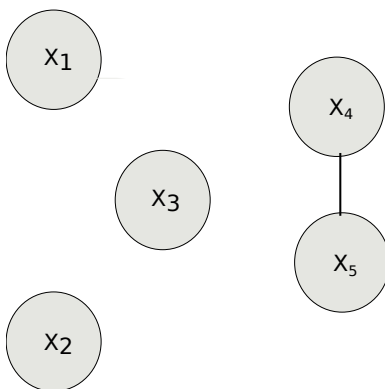
Figure 13: True factor graph representation of the 5-dimensional joint density $p(X_1, \ldots, X_5) = f_\mu(X_1, 1) f_\mu(X_2, 1) f_\mu(X_3, 1) f_\mu(X_4, X_5, 2)$.

For the surrogate cross entropy test, we set $N = .2 * 10^4$, $M = .8 * 10^4$ and $k = 20$. We note that the maximum value of $m_1 + m_2$ for the above set of tests is 8 and that $V/8 > (N + M)$. This choice of $N$ and $M$ therefore ensures that there are enough samples $V$ to guarantee sufficient number of independent samples for estimating individual entropies (see Section 5).

The table below lists the probability (experimental/theoretical prediction[1]) of choosing the false model over the true model for the various tests.

| Same true vs Same false | $l_T$ vs $l_F$ | $m_T$ vs $m_F$ | $n_T$ vs $n_F$ |
|---|---|---|---|
| Error (Exp/Theor) | 0.071/0.032 | 0.067/0.066 | 0.068/0.028 |
| High true vs Low false | $l_T$ vs $m_F$ | $l_T$ vs $n_F$ | $m_T$ vs $n_F$ |
| Error (Exp/Theor) | 0/0 | 0/0 | 0/0 |
| Low true vs High false | $m_T$ vs $l_F$ | $n_T$ vs $l_F$ | $n_T$ vs $m_F$ |
| Error  (Exp/Theor) | 0.689/0.732 | 0.995/1.000 | 0.691/0.665 |

**Explanation** For the class of models above, the set of constants $\{c_{n_i 1}, c_{l_j 1}\}$ are always negative. As a result, when comparing a high dimensional model to a low dimensional model, the additional bias will strongly tilt the test statistic towards the higher dimensional model. As a result, there is a greater chance of detecting the higher dimension model in the surrogate cross entropy test, irrespective of whether the higher dimensional model is true or false.

To elaborate, when the high dimensional model is true and the low dimensional model is false, the bias will further tilt the test statistic towards the high dimensional model, resulting in zero false detections. On the other hand, when the low dimensional model is true, the bias in the surrogate test statistic deviates towards the high dimensional model, resulting in

---

[1]The theoretical prediction requires estimation of constants $c_{l_i 1}, c_{l_i 2}$ and $c_{l_i 3}$. These constants were estimated from the data using oracle Monte Carlo methods which utilized the true form of the density $p$. In practice, when the true form of $p$ is never known, we adopt methods given by [38] to estimate these constants from data.

a high number of false detections. When we compare factor graph models within the same class of dimension, the bias from the cross entropy estimates for each model nearly cancel, resulting in a surrogate test statistic with much smaller bias as compared to the above two cases. As a result, the number of false detections is correspondingly low when comparing models within the same dimension.

By the same argument, for factor graph models where the set of constants $\{c_{n_i 1}, c_{l_j 1}\}$ are positive, we can conclude that the surrogate test statistic will be biased towards lower dimensional models.

# 9 Application to intrinsic dimension estimation

In this work we introduce a new dimensionality estimator that is based on fluctuations of the sizes of nearest neighbor balls centered at a subset of the data points. In this respect it is similar to Costa's $k$-nearest neighbor (kNN) graph dimension estimator [9] and to Farahmand's dimension estimator based on nearest neighbor distances [14]. The estimator can also be related to the Leonenko's Rényi entropy estimator [27]. However, unlike these estimators, our new dimension estimator is derived directly from a mean squared error (M.S.E.) optimality condition for partitioned kNN estimators of multivariate density functionals. This guarantees that our estimator has the best possible M.S.E. convergence rate among estimators in its class. Empirical experiments are presented that show that this asymptotic optimality translates into improved performance in the finite sample regime.

## 9.1 Problem formulation

Let $\mathcal{Y} = \{\mathbf{Y}_1, \ldots, \mathbf{Y}_T\}$ be $T$ independent and identically distributed sample realizations in $\mathbb{R}^D$ distributed according to density $f$. Assume the random vectors in $\mathcal{Y}$ are constrained to lie on a d-dimensional Riemannian submanifold $\mathcal{S}$ of $\mathbb{R}^D$ $(d < D)$. We are interested in estimating the intrinsic dimension $d$.

## 9.2 Log-length statistics

Let $\gamma > 0$ be any arbitrary number and $\alpha = \gamma/d$. Partition the $T$ samples in $\mathcal{Y}$ into two disjoint sets $\mathcal{X}$ and $\mathcal{Z}$ of size $\lfloor T/2 \rfloor$ each. Denote the samples of $\mathcal{X}$ as $\mathcal{X} = \{\mathbf{X}_1, \ldots, \mathbf{X}_{\lfloor T/2 \rfloor}\}$ and $\mathcal{Z}$ as $\mathcal{Z} = \{\mathbf{Z}_1, \ldots, \mathbf{Z}_{\lfloor T/2 \rfloor}\}$.

Partition $\mathcal{X}$ into $N$ 'target' and $M$ 'reference' samples $\{\mathbf{X}_1, \ldots, \mathbf{X}_N\}$ and $\{\mathbf{X}_{N+1}, \ldots, \mathbf{X}_{\lfloor T/2 \rfloor}\}$ respectively with $N + M = \lfloor T/2 \rfloor$. Partition $\mathcal{Z}$ in an identical manner. Now consider the following statistics based on the partitioning of sample space:

$$\mathbf{L_k}(\mathcal{X}) = \frac{\gamma}{N} \sum_{i=1}^{N} \log\left(\mathbf{R_k}(\mathbf{X}_i)\right),$$

30

where $\mathbf{R_k}(\mathbf{X}_i)$ is the Euclidean $k$ nearest neighbor ($k$NN) distance from the target sample $\mathbf{X}_i$ to the $M$ reference samples $\{\mathbf{X}_{N+1}, \ldots, \mathbf{X}_{\lfloor T/2 \rfloor}\}$ . This partitioning of samples is illustrated in Fig. 14.
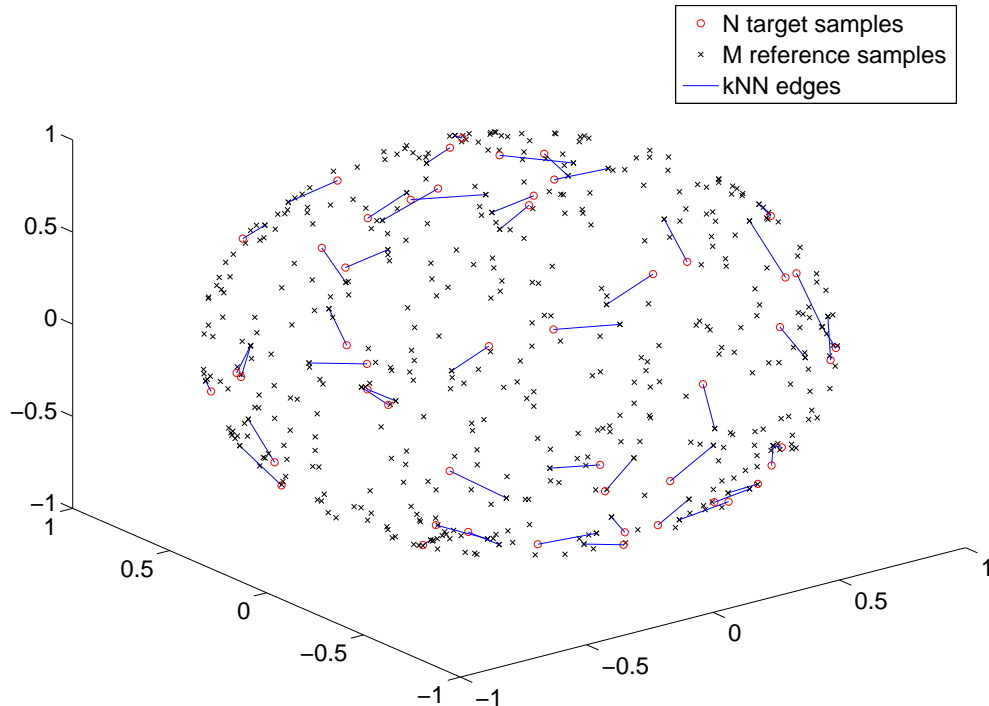


Figure 14: kNN edges on sphere manifold with uniform distribution for $d = 2$, $D = 3$, and $k = 5$.

## 9.3 Relation to $k$NN density estimates

Under the condition that $k/M$ is small, the Euclidean $k$NN distance $\mathbf{R_k}(\mathbf{X}_i)$ approximates the $k$NN distance on the submanifold $\mathcal{S}$. The $k$NN density estimate [31] of $f$ at $\mathbf{X}_i$ based on the $M$ samples $\mathbf{X}_{N+1}, \ldots, \mathbf{X}_{N+M}$ is then given by

$$\hat{\mathbf{f}}_\mathbf{k}(\mathbf{X_i}) = \frac{k-1}{M} \frac{1}{c_d \mathbf{R_k}(\mathbf{X}_i)^d} = \frac{k-1}{M} \frac{1}{\mathbf{V_k}(\mathbf{X}_i)},$$

where $c_d$ is the volume of the unit ball in $d$ dimensions and therefore $\mathbf{V_k(X_i)}$ is the volume of the $k$NN ball. This implies that $\mathbf{L_k}(\mathcal{X})$ can be rewritten as follows:

$$
\begin{aligned}
\mathbf{L_k}(\mathcal{X}) &= \frac{\gamma}{N} \sum_{i=1}^{N} \log\left(\mathbf{R_k(X_i)}\right) \\
&= \log\left(\frac{k-1}{Mc_d}\right)^{\alpha} + \frac{1}{N}\sum_{i=1}^{N} \log\left(\hat{\mathbf{f}}_\mathbf{k}(\mathbf{X}_i)\right)^{-\alpha} \\
&= \alpha \log(k-1) - \frac{\alpha}{N}\sum_{i=1}^{N} \log\hat{\mathbf{f}}_\mathbf{k}(\mathbf{X}_i) \\
&\quad - \alpha \log(c_d M).
\end{aligned}
\tag{29}
$$

As eq. (29) indicates, the log-length statistics is linear with respect to $\log(k-1)$ with a slope of $\alpha$. This prompts the idea of estimating $\alpha$ (and later $d$) from the slope of $\mathbf{L_k}(\mathcal{X})$ as a function of $\log(k-1)$.

## 9.4  Intrinsic dimension estimate based on varying bandwidth $k$

Let $k_1$ and $k_2$ be two different choices of bandwidth parameters. Let $\mathbf{L_{k_1}}(\mathcal{X})$ and $\mathbf{L_{k_2}}(\mathcal{Z})$ be the length statistics evaluated at bandwidths $k_1$ and $k_2$ using data $\mathcal{X}$ and $\mathcal{Z}$ respectively. A natural choice for the estimate of $\alpha$ would then be

$$
\begin{aligned}
\hat{\alpha} &= \frac{\mathbf{L_{k_2}}(\mathcal{Z}) - \mathbf{L_{k_1}}(\mathcal{X})}{\log(k_2-1) - \log(k_1-1)} \\
&= \alpha + \frac{\nu}{N}\sum_{i=1}^{N}\left(\log\hat{\mathbf{f}}_\mathbf{k_2}(\mathbf{Z}_i) - \log\hat{\mathbf{f}}_\mathbf{k_1}(\mathbf{X}_i)\right) \\
&= \alpha + \nu(\hat{\mathbf{E}}_\mathbf{k_2}(\mathcal{Z}) - \hat{\mathbf{E}}_\mathbf{k_1}(\mathcal{X})),
\end{aligned}
$$

where

$$
\hat{\mathbf{E}}_\mathbf{k}(\mathcal{X}) = \frac{1}{N}\sum_{i=1}^{N}\log(\hat{\mathbf{f}}_\mathbf{k}(\mathbf{X_i})),
$$

and $\nu = -\alpha/\log((k_2-1)/(k_1-1))$. The intrinsic dimension estimate is related to $\hat{\alpha}$ by the simple relation $\hat{\mathbf{d}} = \gamma/\hat{\alpha}$.

## 9.5 Statistical properties of intrinsic dimension estimate

We can relate the error in estimation of $\alpha$ to the error in dimension estimation as follows:

$$
\begin{aligned}
\hat{\mathbf{d}} - d &= \gamma \left( \frac{1}{\hat{\alpha}} - \frac{1}{\alpha} \right) \\
&= \gamma \frac{\alpha - \hat{\alpha}}{\hat{\alpha}\alpha} \\
&= -\frac{\gamma}{\alpha^2}(\hat{\alpha} - \alpha) + o(\hat{\alpha} - \alpha).
\end{aligned}
$$

Define $\kappa = -\gamma\nu/\alpha^2$. We recognize that the density functional estimate $\hat{\mathbf{E}}_{\mathbf{k}}(\mathcal{X})$ is in the form of the plug-in estimators introduced in this report. Using the results on the bias, variance and asymptotic distribution of the density functional estimate $\hat{\mathbf{E}}_{\mathbf{k}}(\mathcal{X})$ established in this report and the above relation between the errors $\hat{\mathbf{d}} - d$ and $\hat{\alpha} - \alpha$, we then have the following statistical properties for the estimate $\hat{\mathbf{d}}$:

**Estimator bias**

$$
\begin{aligned}
\mathbb{E}[\hat{\mathbf{d}}] - d &= \kappa c_{b_1} \left( \left( \frac{k_2}{M} \right)^{2/d} - \left( \frac{k_1}{M} \right)^{2/d} \right) \\
&+ \kappa c_{b_2} \left( \left( \frac{1}{k_2} \right) - \left( \frac{1}{k_1} \right) \right) \\
&+ o \left( \frac{1}{k_1} + \frac{1}{k_2} + \left( \frac{k_1}{M} \right)^{2/d} + \left( \frac{k_2}{M} \right)^{2/d} \right).
\end{aligned}
$$

**Estimator variance**

$$
\mathbb{V}(\hat{\mathbf{d}}) = 2\kappa^2 c_v \left( \frac{1}{N} \right) + o \left( \frac{1}{M} + \frac{1}{N} \right).
$$

**Central limit theorem**

Let $\mathbf{Z}$ be a standard normal random variable. Then,

$$
\lim_{N,M \to \infty} Pr \left( \frac{\hat{\mathbf{d}} - \mathbb{E}[\hat{\mathbf{d}}]}{\sqrt{2\kappa^2 c_v/N}} \leq \alpha \right) = Pr(\mathbf{Z} \leq \alpha).
$$

## 9.6 Optimal selection of parameters

We have theoretical expressions for the mean square error (M.S.E) of the dimension estimate $\hat{\mathbf{d}}$, which we can optimize over the free parameters $k_1$, $k_2$, $N$ and $M$. We restrict our attention

to the case $k_2 = 2k$; $k_1 = k$. The M.S.E. of $\hat{\mathbf{d}}$ (ignoring higher order terms) is given by

$$
\begin{aligned}
\text{M.S.E.}(\hat{\mathbf{d}}) &= (\mathbb{E}[\hat{\mathbf{d}}] - d)^2 + \mathbb{V}[\hat{\mathbf{d}}] \\
&= \left( C_{b_1} \left( \frac{k}{M} \right)^{2/d} + C_{b_2} \left( \frac{1}{k} \right) \right)^2 \\
&\quad + C_v \left( \frac{1}{N} \right).
\end{aligned} \tag{30}
$$

where $C_{b_1} = \kappa 2^{(2/d - 1)}$, $C_{b_2} = \kappa/4$ and $C_v = 2\kappa^2 c_v$.

## Optimal choice of bandwidth

The optimal value of $k$ w.r.t the M.S.E. is given by

$$
k_{opt} = \lfloor k_0 M^{\frac{2}{2+d}} \rfloor. \tag{31}
$$

where the constant $k_0 = (|C_{b_2}| d/2 |C_{b_1}|)^{\frac{d}{d+2}}$.

## Optimal partitioning of sample space

Under the constraint that $N + M = \lfloor T/2 \rfloor$ is fixed, the optimal choice of $N$ as a function of $M$ is then given by

$$
N_{opt} = \lfloor N_0 M^{\frac{6+d}{2(2+d)}} \rfloor, \tag{32}
$$

where the constant $N_0 = \frac{\sqrt{C_v(2+d)}}{2b_0}$.

## 9.7  Improved estimator based on correlated error

Consider the following alternative estimator for $\alpha$:

$$
\begin{aligned}
\tilde{\alpha} &= \frac{\mathbf{L}_{\mathbf{k_2}}(\mathcal{X}) - \mathbf{L}_{\mathbf{k_1}}(\mathcal{X})}{\log(k_2 - 1) - \log(k_1 - 1)} \\
&= \alpha + \kappa(\hat{\mathbf{E}}_{\mathbf{k_2}}(\mathcal{X}) - \hat{\mathbf{E}}_{\mathbf{k_1}}(\mathcal{X})),
\end{aligned}
$$

and the corresponding density estimate $\tilde{\mathbf{d}}$ which satisfies

$$
\tilde{\mathbf{d}} - d = -\frac{\gamma}{\alpha^2}(\tilde{\alpha} - \alpha) + o(\hat{\alpha} - \alpha),
$$

where both the length statistics at bandwidths $k_1$ and $k_2$ are evaluated using the same sample $X$. The density functional estimates $\hat{\mathbf{E}}_{\mathbf{k_1}}(\mathcal{X})$ and $\hat{\mathbf{E}}_{\mathbf{k_2}}(\mathcal{X})$ will be highly correlated (as compared to the independent quantities $\hat{\mathbf{E}}_{\mathbf{k_1}}(\mathcal{X})$ and $\hat{\mathbf{E}}_{\mathbf{k_2}}(\mathcal{Z})$). This implies that the variance
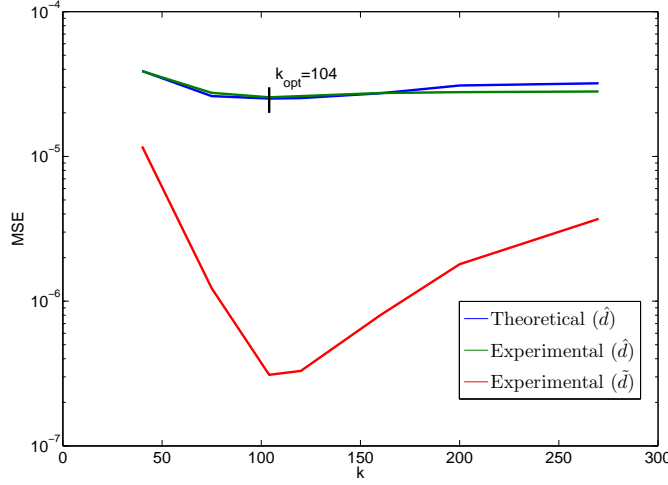
Figure 15: Comparison of theoretically predicted and experimental M.S.E. for varying choices of $k$. The experimental performance of the estimator $\hat{\mathbf{d}}$ is in excellent agreement with the theoretical expression and, as predicted by our theory, the modified estimator $\tilde{\mathbf{d}}$ significantly outperforms $\hat{\mathbf{d}}$.

of the difference $\hat{\mathbf{E}}_{\mathbf{k_2}}(\mathcal{X}) - \hat{\mathbf{E}}_{\mathbf{k_1}}(\mathcal{X})$ will be smaller when compared to $\hat{\mathbf{E}}_{\mathbf{k_2}}(\mathcal{Z}) - \hat{\mathbf{E}}_{\mathbf{k_1}}(\mathcal{X})$, (while the expectation remains the same).

Since the estimator bias is unaffected by this modification, the variance reduction suggests that $\tilde{d}$ will be an improved estimator as compared to $\hat{\mathbf{d}}$ in terms of M.S.E.. In order to obtain statistical properties for the improved estimator $\tilde{\mathbf{d}}$ (equivalent to the properties developed in Section 9.5 for the original estimator $\hat{\mathbf{d}}$), we need to analyze the joint distribution between $\hat{\mathbf{f}}_{\mathbf{k_1}}(X_i)$ and $\hat{\mathbf{f}}_{\mathbf{k_2}}(X_j)$ for two distinct values $k_1$ and $k_2$. Our theory, at present, cannot address the case of distinct bandwidths $k_1$ and $k_2$.

Since the estimate $\tilde{\mathbf{d}}$ has smaller M.S.E. compared to $\hat{\mathbf{d}}$, M.S.E. predictions for the estimate $\hat{\mathbf{d}}$ can serve as upper bounds on the M.S.E. performance of the improved estimate $\tilde{\mathbf{d}}$.

## 9.8 Simulations

We generate $T = 10^5$ samples $\mathcal{B}$ drawn from a $d = 2$ mixture density $f_m = .8 f_\beta + .2 f_u$, where $f_\beta$ is the product of two 1 dimensional marginal beta distributions with parameters $\alpha = 2$, $\beta = 2$ and $f_u$ is a uniform density in 2 dimensions. These samples are then projected to a 3-dimensional hyperplane in $\mathbb{R}^3$ by applying the transformation $\mathcal{Y} = U\mathcal{B}$ where $U$ is a $3 \times 2$ random matrix whose columns are orthonormal. We apply our intrinsic dimension estimates on the samples $\mathcal{Y}$.
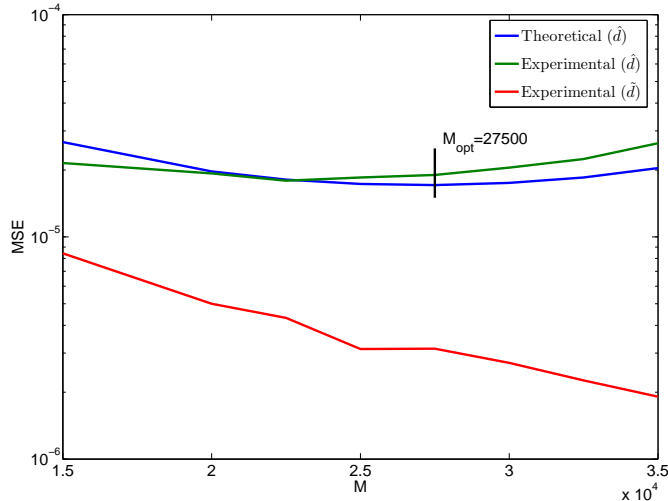
35

Figure 16: Comparison of theoretically predicted and experimental M.S.E. for varying choices of $M$. The experimental performance of the estimator $\hat{\mathbf{d}}$ is in excellent agreement with the theoretical expression and, as predicted by our theory, the modified estimator $\tilde{\mathbf{d}}$ significantly outperforms $\hat{\mathbf{d}}$.

## Optimal selection of free parameters

In our first experiment, we theoretically compute the optimal choice of $k$ for a fixed partition with $M = 3.5 \times 10^4$ and $N = 1.5 \times 10^4$. We then show the variation of the theoretical and experimental M.S.E. of the estimate $\hat{\mathbf{d}}$ and the experimental M.S.E. of the improved estimate $\tilde{\mathbf{d}}$ with changing bandwidth $k$ in Fig. 15. In our second experiment, we compute the optimal partition according to eq. (32) and show the variation of M.S.E. with varying choices of partition in Fig. 16.

From our experiments, we see that there is good agreement between our theory and simulations. As a consequence, we find the theoretically predicted optimal choices of $k, N$ and $M$ to minimize the observed M.S.E.. In addition, as predicted by our theory, the modified estimator $\tilde{\mathbf{d}}$ significantly outperforms $\hat{\mathbf{d}}$. The theoretically predicted M.S.E. for $\hat{\mathbf{d}}$ therefore serves as a strict upper bound for the M.S.E. of the improved estimator $\tilde{\mathbf{d}}$.

## Comparison of dimension estimation methods

We compare the performance of our proposed dimension estimators to the estimated proposed by Frahmand et. al. [14] (denote as $\hat{\mathbf{d}}_f$) and Costa et. al. [9] (denote as $\hat{\mathbf{d}}_j$).

Expressions for the optimal bandwidth $k$ (eq. (4)) and partition $N, M$ (eq. (32)) depend on the unknown intrinsic dimension $d$ and constants $c_{b_1}$, $c_{b_2}$ and $c_v$ which depend on unknown
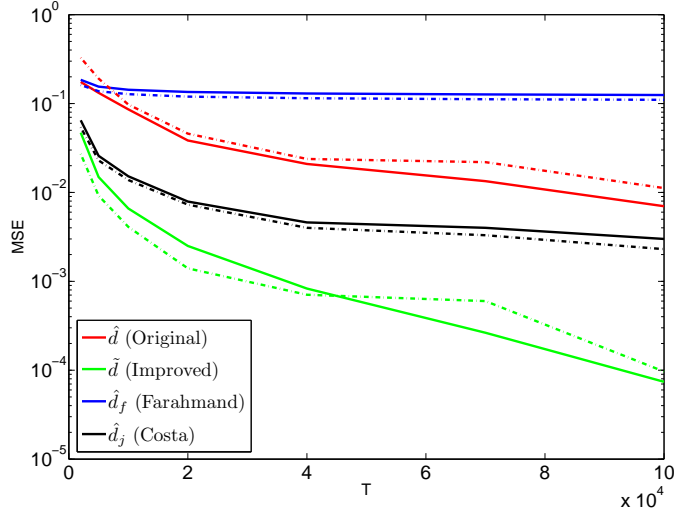
Figure 17: Comparison of performance of dimension estimates (Solid line: Optimal (optimal choice of $k$,$N$ and $M$ as per eq. (4) and eq. (32)); Dashed line: Suboptimal (fixed $k = 20$, $N = T/50$, $M = \lfloor T/2 \rfloor - N$)): The proposed improved kNN distance estimator outperforms all other estimators considered.

density $f$. The constants $c_{b_1}$, $c_{b_2}$ and $c_v$ can be estimated from the data using plug-in methods similar to the ones used by Raykar et. al. [38] for optimal bandwidth selection for kernel density estimation . To establish the potential advantages of our dimension estimators we compare an omniscient optimal form of our estimator, for which the true values of these constants are known, to a suboptimal form of our estimator that does not know the constants.

For the optimal estimator, we theoretically compute the optimal choice for $k$, $N$ and $M$ for different choices of total sample size $T$ (sub-sampled from the initial $10^5$ samples), and use these optimal parameters for the estimators $\hat{\mathbf{d}}$ and $\tilde{\mathbf{d}}$. We use this optimal choice of bandwidth $k$ for the estimators $\hat{\mathbf{d}}_f$ and $\hat{\mathbf{d}}_j$ as well (partitioning not applicable). For the suboptimal estimator, we arbitrarily choose the parameters as follows: fixed $k = 20$, $N = T/50$, $M = \lfloor T/2 \rfloor - N$.

The performance of these estimators as a function of sample size $T$ is shown in Fig. 17. Estimators with optimal choice of parameters are indicated in solid line, and the suboptimal estimators are indicated in dashed lines.

From our experiments we see that the performance of the original estimator $\hat{\mathbf{d}}$ with sub-optimal choice of parameters is marginally inferior when compared to the estimator with optimal choice of parameters. This does not hold for the other estimators as can be expected since the parameters are optimized w.r.t. the performance of $\hat{\mathbf{d}}$.

We note that the improved estimator $\tilde{\mathbf{d}}$ outperforms all other estimators while the performance of our original estimator $\hat{\mathbf{d}}$ is sandwiched between $\hat{\mathbf{d}}_f$ and $\hat{\mathbf{d}}_j$. We conjecture that the
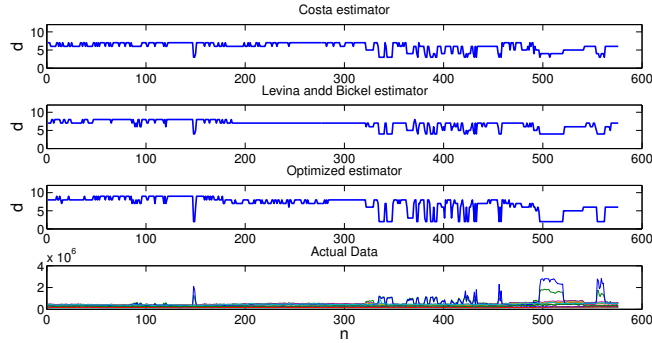
Figure 18: Comparison of performance of dimension estimates for anomaly detection in Abilene network data.

performance of $\hat{\mathbf{d}}_j$ is superior to $\hat{\mathbf{d}}$ for the same reason that $\tilde{\mathbf{d}}$ outperforms $\hat{\mathbf{d}}$: correlated error between different length statistics.

## Anomaly detection in Abilene network data

Anomalies can be detected in router netowrks by estimating the local dimension at each time point and monitoring change in dimension. The data used is the number of packets sent by each of the 11 routers on the abiline network between January 1-2, 2005. A sample is taken every 5 minutes, leading to 576 samples with an extrinsic dimension pf 11.

The performance of different dimension estimators is shown in Fig. 18. We know that simulataneous peaks in router traffic should imply strong correlation between the routers and therefore lower intrinsic dimension. This behaviour is clearly reflected better by the optimized estimator as compared to the estimator of Costa et. al. [9] and Levina and Bickel [28].

# 10   Conclusion

A new class of boundary compensated bipartite k-NN density plug-in estimators was proposed for estimation of smooth non-linear functionals of densities that are strictly bounded strictly away from 0 on their finite support. These estimators, called bipartite plug-in (BPI) estimators, correct for bias due to boundary effects and outperform previous $k$-NN entropy estimators in terms of MSE convergence rate. Expressions for asymptotic bias and variance of the estimator were derived estimator in terms of the sample size, the dimension of the samples and the underlying probability distribution. In addition, a central limit theorem was developed for the proposed BPI estimators. The accuracy of these asymptotic results were validated through simulation and it was established that the theory can be used to specify optimal finite sample estimator tuning parameters such as bandwidth and optimal

partitioning of data samples.

Our theory has two important by-products: (1) We established similarity between the moments of $k$-NN density estimates and kernel density estimates. This in turn implies that plug-in estimators based on $k$-NN density estimators and kernel density estimators have asymptotically equal rates of convergence. (2) We developed an algorithm for detection and correction of density estimates at boundary points for densities with finite support. This correction helps reduce the bias of density estimates at the boundaries of the support of the density, thereby reducing the overall bias of the plug-in estimators.

Using the theory presented in the paper, one can tune the parameters of the plug-in estimator to achieve minimum asymptotic estimation MSE. Furthermore, the theory can be used to specify the minimum necessary sample size required to obtain requisite accuracy. This in turn can be used to predict and optimize performance in applications like structure discovery in graphical models and dimension estimation for support sets of low intrinsic dimension. We applied our theory to the problem of estimating Shannon entropy and Shannon mutual information. Furthermore, we used the Shannon entropy estimator to discover structure in high dimensional data and to determine the intrinsic dimension of data samples.

For the reader's convenience, the notation used in this paper is listed in the table below.

| Notation | Description |
|---|---|
| $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ | BPI estimator (1) |
| $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)$ | BPI estimator with bias compensation (9) |
| $g_1(k,M), g_2(k,M)$ | Bias correction factors |
| $\mathcal{S}$ | Support of density $f$ |
| $d$ | dimension of support $\mathcal{S}$ |
| $c_d$ | unit ball volume in $d$ dimensions |
| $\{\mathbf{X}_1,\ldots,\mathbf{X}_T,\mathbf{Y},\mathbf{Z}\}$ | $T+2$ independent realizations drawn from $f$ |
| $\mathcal{X}_N$ | $\{\mathbf{X}_1,\ldots,\mathbf{X}_N\}$ |
| $\mathcal{X}_M$ | $\{\mathbf{X}_{N+1},\ldots,\mathbf{X}_{N+M}\}$ |
| $\mathcal{S}_I$ | Interior of support |
| $\mathcal{I}_N$ | Interior points subset of $\mathcal{X}_N$ |
| $\mathcal{B}_N$ | Boundary points subset of $\mathcal{X}_N$ |
| $\mathbf{Z}_{-1}$ | Closest interior point to $\mathbf{Z}$; $\mathbf{Z}_{-1} = \arg\min_{x \in \mathcal{S}_I} d(x, \mathbf{Z})$ |
| $\mathbf{X}_{n(i)}$ | $\mathbf{X}_{n(i)} \in \mathcal{I}_N$ is the interior sample point that is closest to $\mathbf{X}_i \in \mathcal{B}_N$ |
| $\delta$ | Constant; $\delta \in (2/3, 1)$ |
| $\epsilon_{BC} = N\exp(-3k^{(1-\delta)})$ | Probability of misclassification of $x \in \mathcal{S} - \mathcal{S}_I$ as interior point |
| $\mathbf{d}_k(X)$ | $k$-NN ball radius |
| $\mathbf{S}_k(X)$ | $k$-NN ball |
| $\mathbf{V}_k(X)$ | $k$-NN ball volume |
| $\mathbf{P}(X)$ | Coverage function |
| $\hat{\mathbf{f}}_k(X)$ | $k$-NN density estimate |
| $\tilde{\mathbf{f}}_k(X)$ | Boundary corrected $k$-NN density estimate |
| $g^{(n)}(x,y)$ | $n$-th derivative of $g(x,y)$ wrt $x$ |
| $\mathbf{p}$ | beta random variable with parameters $k, M-k+1$ |
| $\alpha_{frac}$ | Proportionality constant; $M = \alpha_{frac}T$ and $N = (1-\alpha_{frac})T$ |
| $\epsilon_0, \epsilon_\infty$ | constants such that $\epsilon_0 \le f(x) \le \epsilon_\infty \ \forall x \in \mathcal{S}$ |
| $2\nu$ | Number of times $f$ is assumed to be differentiable |
| $\lambda$ | Number of times $g(x,y)$ is assumed to be differentiable wrt $x$ |
| $c_1,..,c_5$ | Constants appearing in Theorems III.1, III.2, III.3 and IV.1, IV.2, IV.3 |
| $\mathcal{C}(k)$ | Function which satisfies the rate of decay condition $\mathcal{C}(k) = O(e^{-3k^{(1-\delta)}})$ |
| $k_M$ | $k_M = (k-1)/M$ |
| $\natural(X)$ | The event $\mathbf{P}(X) > (1-p_k)k_M$ |
| $\natural_{-1}(X)$ | The event $\mathbf{P}(X) < (1+p_k)k_M$ |
| $\natural\natural(X)$ | The event $(1-p_k)k_M < \mathbf{P}(X) < (1+p_k)k_M$ |
| $\mathbf{e}_k(X)$ | Error function $\mathbf{e}_k(X) = \hat{\mathbf{f}}_k(X) - \mathbb{E}[\hat{\mathbf{f}}_k(X) \mid X]$ |
| $\mathbf{e}(X)$ | Error function $\mathbf{e}(X) = \tilde{\mathbf{f}}_k(X) - \mathbb{E}[\tilde{\mathbf{f}}_k(X) \mid X]$ |

# Appendices

## A  Uniform kernel density estimation

Throughout this section, we will derive results on moments of the uniform kernel density estimates for points in the set $\mathcal{S}' = \{X : \mathbf{S_u}(X) \subset \mathcal{S}\}$. This definition implies that the density $f$ has continuous partial derivatives of order $2r$ in the uniform ball neighborhood for each $X \in \mathcal{S}'$ where $r$ satisfies the condition $2r(1-t)/d > 1$. This excludes the set of points close to the boundary of the support, where the continuity assumption of the density is not satisfied. We will deal with these points in Appendix C.

Let $\mathbf{X}_1, .., \mathbf{X}_M$ denote $M$ i.i.d realizations of the density f. We will assume that $f$ is continuously differentiable evrywhere in the interior of the sWe seek to estimate the density at $X$ from the $M$ i.i.d realizations $\mathbf{X}_1, .., \mathbf{X}_M$. Let $c_d$ denote the volume of a unit hyper-sphere in $d$ dimensions. The uniform kernel density estimator is defined as follows:

### A.1  Uniform kernel density estimator

The *uniform kernel* density estimator is defined below. The volume of the uniform kernel is given by

$$V_u(X) = \frac{k}{M}, \tag{33}$$

and the kernel region is given by

$$S_u(X) = \{Y : c_d||X - Y||^d \leq V_u\}. \tag{34}$$

$\mathbf{l_u}(X)$ denotes the number of points falling in $S_u(X)$

$$\mathbf{l_u}(X) = \Sigma_{i=1}^{M} 1_{X_i \in S_u(X)}, \tag{35}$$

and the *uniform kernel* density estimator is defined by

$$\hat{\mathbf{f}}_{\mathbf{u}}(X) = \frac{\mathbf{l_u}(X)}{MV_u(X)}. \tag{36}$$

The *coverage* of the *uniform kernel* is defined as

$$U(X) = \int_{S_u(X)} f(z)dz = \mathbb{E}[1_{\mathbf{Z} \in S_u(X)}]. \tag{37}$$

We observe that $\mathbf{l_u}(X)$ is a binomial random variable with parameters $M$ and $U(X)$. Figure 19 illustrates the *uniform kernel* density estimate.
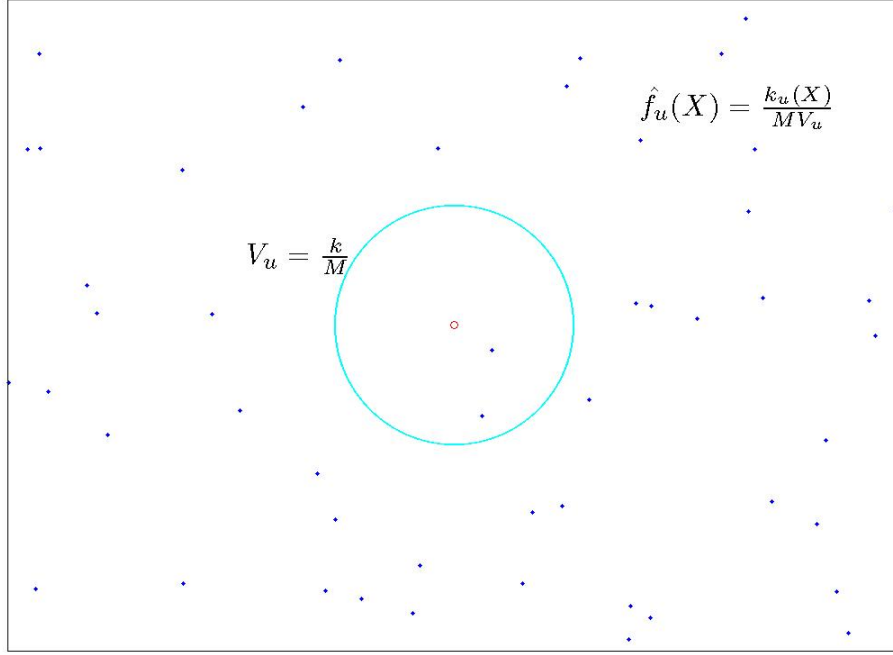


$$\hat{f}_u(X) = \frac{k_u(X)}{MV_u}$$

$$V_u = \frac{k}{M}$$

Figure 19: Uniform kernel density estimator.

## A.2 Taylor series expansion of coverage

**We assume that the density $f$ has continuous partial derivatives of third order** in a neighborhood of $X$. For small volumes $V_u(X)$ (which is equivalent to the condition that $k/M$ is small), we can represent the coverage function $U(X)$ by using a third order Taylor series expansion of $f$ about about $X$ [31].

$$
\begin{aligned}
U(X) &= \int_{S_u(X)} f(Z)dZ \\
&= f(X)V_u(X) + c(X)V_u^{1+2/d}(X) + o(V_u^{1+2/d}(X)) \\
&= f(X)\frac{k}{M} + c(X)\left(\frac{k}{M}\right)^{1+2/d} + o\left(\left(\frac{k}{M}\right)^{1+2/d}\right),
\end{aligned}
\tag{38}
$$

where $c(X) = \Gamma^{(2/d)}(\frac{n+2}{2})tr[\nabla^2(f(X))]$.

## A.3 Concentration inequalities for uniform kernel density

Because $\mathbf{l_u}(X)$ is a binomial random variable, we can apply standard Chernoff inequalities to obtain concentration bounds on the density estimate. $\mathbf{l_u}(X)$ is a binomial random variable with parameters $M$ and $U(X)$.

### A.3.1 Concentration around true density

For $0 < p < 1/2$,
$$Pr(\mathbf{l_u}(X) > (1+p)MU(X)) \leq e^{-MU(X)p^2/4}, \tag{39}$$

and
$$Pr(\mathbf{l_u}(X) < (1-p)MU(X)) \leq e^{-MU(X)p^2/4}. \tag{40}$$

Using the Taylor expansion of coverage, we then have
$$Pr(\hat{\mathbf{f}}_\mathbf{u}(X) > (1+p)(f(X) + O((k/M)^{2/d}))) \leq\sim e^{-p^2kf(X)/4}, \tag{41}$$

and
$$Pr(\hat{\mathbf{f}}_\mathbf{u}(X) < (1-p)(f(X) + O((k/M)^{2/d}))) \leq\sim e^{-p^2kf(X)/4}. \tag{42}$$

This then implies that
$$Pr(\hat{\mathbf{f}}_\mathbf{u}(X) > (1+p)f(X)) \leq\sim e^{-p^2kf(X)/4}, \tag{43}$$

and
$$Pr(\hat{\mathbf{f}}_\mathbf{u}(X) < (1-p)f(X)) \leq\sim e^{-p^2kf(X)/4}. \tag{44}$$

Let $\mathbf{X}$ be a random variable with density $f$ independent of the $M$ i.i.d realizations $\mathbf{X}_1, .., \mathbf{X}_M$. Then,

$$\begin{aligned}
Pr(\hat{\mathbf{f}}_\mathbf{u}(\mathbf{X}) > (1+p)f(\mathbf{X})) &= \mathbb{E}_\mathbf{X}[Pr(\hat{\mathbf{f}}_\mathbf{u}(\mathbf{X}) > (1+p)f(\mathbf{X}))] \\
&\leq \mathbb{E}[\sim (e^{-p^2kf(\mathbf{X})/4})] \\
&= \sim e^{-p^2k/4},
\end{aligned} \tag{45}$$

and

$$\begin{aligned}
Pr(\hat{\mathbf{f}}_\mathbf{u}(\mathbf{X}) < (1-p)f(\mathbf{X})) &= \mathbb{E}_\mathbf{X}[Pr(\hat{\mathbf{f}}_\mathbf{u}(\mathbf{X}) < (1-p)f(\mathbf{X}))] \\
&\leq \mathbb{E}[\sim (e^{-p^2kf(\mathbf{X})/4})] \\
&= \sim e^{-p^2k/4}.
\end{aligned} \tag{46}$$

### A.3.2 Concentration away from $0$

We can also bound the density estimate away from 0 as follows:

$$
\begin{aligned}
Pr(\hat{\mathbf{f}}_{\mathbf{u}}(\mathbf{X}) = 0) &= \mathbb{E}_{\mathbf{X}}[Pr(\hat{\mathbf{f}}_{\mathbf{u}}(\mathbf{X}) = 0] \\
&= \mathbb{E}[(1 - U(X))^M] \\
&= \mathbb{E}[(1 - (kf(X) + o(k)/M)^M] \\
&= \mathbb{E}[((1 - (kf(X) + o(k)/M)^{M/(kf(X)+o(k))})^{kf(X)+o(k)}] \\
&= \mathbb{E}[\sim (1/e)^{kf(X)+o(k)}] \\
&= \sim e^{-k}.
\end{aligned}
\tag{47}
$$

## A.4 Central Moments

Define the error function of the uniform kernel density,

$$
\mathbf{e}_{\mathbf{u}}(X) = \hat{\mathbf{f}}_{\mathbf{u}}(X) - \mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(X)].
\tag{48}
$$

The probability mass function of the binomial random variable $\mathbf{l}_{\mathbf{u}}(X)$ is given by

$$
Pr(\mathbf{l}_{\mathbf{u}}(X) = l_x) = \binom{M}{l_x}(U(X))^{l_x}(1 - U(X))^{M-l_x}.
$$

Since $\mathbf{l}_{\mathbf{u}}(X)$ is a binomial random variable, we can easily obtain moments of the uniform kernel density estimate. These are listed below.

First Moment:

$$
\begin{aligned}
\mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(X)] - f(X) &= \frac{M}{k}U(X) - f(X) \\
&= c(X)\left(\frac{k}{M}\right)^{2/d} + o\left(\left(\frac{k}{M}\right)^{2/d}\right).
\end{aligned}
\tag{49}
$$

Second Moment:

$$
\begin{aligned}
\mathbb{V}[\hat{\mathbf{f}}_{\mathbf{u}}(X)] &= \mathbb{E}[\mathbf{e}_{\mathbf{u}}^2(X)] \\
&= \frac{M}{k^2}U(X)(1 - U(X)) \\
&= f(X)\frac{1}{k} + o\left(\frac{1}{k}\right).
\end{aligned}
\tag{50}
$$

Higher Moments: For any integer $r \geq 3$,

$$
\mathbb{E}[\mathbf{e}_u^r(X)] = O\left(\frac{1}{k^{r/2}}\right).
\tag{51}
$$

## A.5 Covariance

Let $X$ and $Y$ be two distinct points. Clearly the density estimates at $X$ and $Y$ are not independent. We expect the density estimates to have positive covariance if $X$ and $Y$ are close and have negative covariance if $X$ and $Y$ are far. This is illustrated in Figure 20.
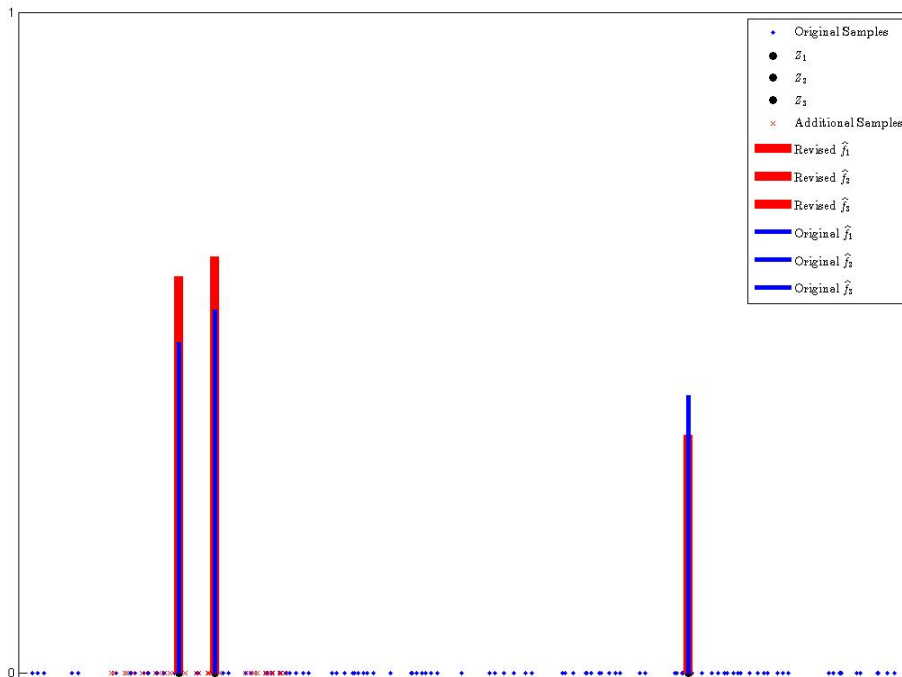


Figure 20: Covariance between uniform kernel density estimates.

Observe that the uniform kernels are disjoint for the set of points given by $\Psi_u := \{X, Y\}$ : $||X - Y|| \geq 2(k/c_d M)^{1/d}$, and have finite intersection on the complement of $\Psi_u$. Indeed we will show that when the uniform balls intersect (and therefore $X$ and $Y$ are close), the density estimates have positive covariance and that they have negative covariance when the uniform kernels are disjoint. Intersecting and disjoint balls are illustrated in Figure 21.

Define,

$$U(X, Y) := \mathbb{E}[1_{\mathbf{Z} \in S_u(X)} 1_{\mathbf{Z} \in S_u(Y)}]. \tag{52}$$

**Intersecting balls**

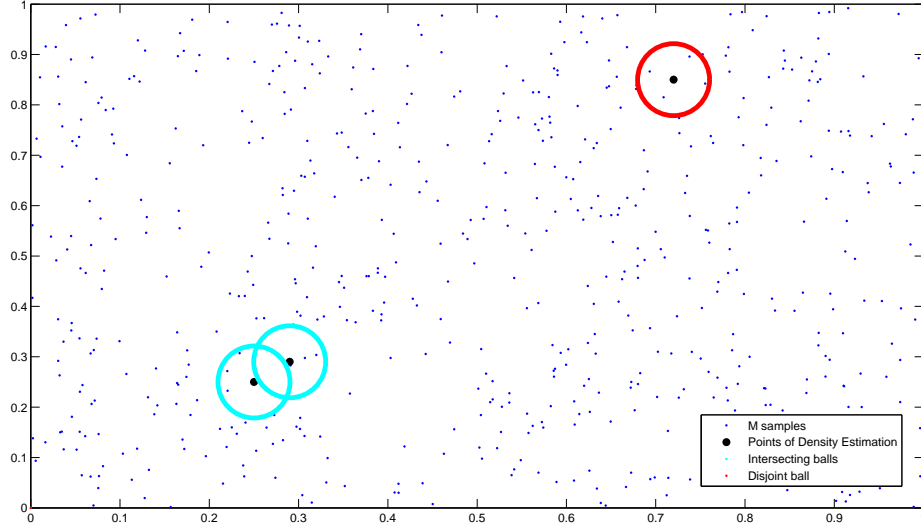**Lemma A.1.** *For a fixed pair of points $\{X, Y\} \in \Psi_u$,*

45

Figure 21: Intersecting and disjoint balls.

$$Cov[\mathbf{e_u}(X), \mathbf{e_u}(Y)] = \frac{-f(X)f(Y)}{M} + o\left(\frac{1}{M}\right).$$

*Proof.* For $\{X, Y\} \in \Psi_u$, we have that $1_{\mathbf{Z} \in S_u(X)} 1_{\mathbf{Z} \in S_u(Y)} = 0$ and therefore $U(X, Y) = 0$. We then have,

$$
\begin{aligned}
Cov[\mathbf{e_u}(X), \mathbf{e_u}(Y)] &= \mathbb{E}[(\hat{\mathbf{f}}_\mathbf{u}(X) - \mathbb{E}[\hat{\mathbf{f}}_\mathbf{u}(X)])(\hat{\mathbf{f}}_\mathbf{u}(Y) - \mathbb{E}[\hat{\mathbf{f}}_\mathbf{u}(Y)])] \\
&= \frac{M}{k^2} \mathbb{E}[(1_{\mathbf{Z} \in S_u(X)} - U(X))(1_{\mathbf{Z} \in S_u(Y)} - U(Y))] \\
&= \frac{M}{k^2} \mathbb{E}[1_{\mathbf{Z} \in S_u(X)} 1_{\mathbf{Z} \in S_u(Y)} - U(X)U(Y)] \\
&= \frac{M}{k^2}(U(X, Y) - U(X)U(Y)) \\
&= -\frac{M}{k^2}[U(X)U(Y)] = \frac{-f(X)f(Y)}{M} + o\left(\frac{1}{M}\right).
\end{aligned}
$$

$\square$

**Disjoint balls** For $\{X, Y\} \in \Psi_u^c$, there is no closed form expression for the covariance. However we have the following lemmas:

Let $R_u(X)$ and $R_u(Y)$ denote the (constant and equal) radii of the uniform balls respectively. Define $\aleph(\|X - Y\|/R_u(X)) = V(S_u(X) \cap S_u(Y))/V_u(X)$ where $V(S_u(X) \cap S_u(Y))$ is the volume of the intersection of the two balls.

We observe that,

$$
\begin{aligned}
\aleph(\|X - Y\|/R_u(X)) &= V(S_u(X) \cap S_u(Y))/V_u(X) \\
&= \frac{V[1_{\mathbf{Z} \in B(0, R_u(X))} 1_{\mathbf{Z} \in B(\|Y - X\|, R_u(Y))}]}{V_u(X)} \\
&= \frac{V[1_{\mathbf{Z} \in B(0,1)} 1_{\mathbf{Z} \in B(\|Y - X\|/R_u(X), 1)}]}{V[1_{\mathbf{Z} \in B(0,1)}]} \\
&= O(1).
\end{aligned}
\tag{53}
$$

Because $f$ is assumed to be continuous, we have

$$
U(X, Y) = \mathbb{E}[1_{\mathbf{Z} \in S_u(X)} 1_{\mathbf{Z} \in S_u(Y)}] = [f(X) + o(1)] V(S_u(X) \cap S_u(Y)).
\tag{54}
$$

**Lemma A.2.** *For a fixed pair of points $\{X, Y\} \in \Psi_u{}^c$,*

$$
Cov[\mathbf{e_u}(X), \mathbf{e_u}(Y)] = O(1/k).
$$

*Proof.*

$$
\begin{aligned}
\frac{M}{k^2} U(X, Y) &= \frac{M}{k^2}[f(X) + o(1)] V(S_u(X) \cap S_u(Y)) \\
&= \frac{f(X) + o(1)}{k} \frac{V(B_X \cap B_Y)}{V_u(X)} \\
&= \frac{f(X) + o(1)}{k} \aleph(\|X - Y\|/R_u(X)) \\
&= \frac{f(X)}{k} \aleph(\|X - Y\|/R_u(X)) + o(1/k) \\
&= O(1/k).
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
Cov[\mathbf{e_u}(X), \mathbf{e_u}(Y)] &= \mathbb{E}[(\hat{\mathbf{f}}_{\mathbf{u}}(X) - \mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(X)])(\hat{\mathbf{f}}_{\mathbf{u}}(Y) - \mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(Y)])] \\
&= \frac{M}{k^2}(U(X, Y) - U(X)U(Y)) \\
&= \frac{M}{k^2}U(X, Y) - \frac{M}{k^2}U(X)U(Y) \\
&= O(1/k) - \Theta(1/M) \\
&= O(1/k).
\end{aligned}
$$

$$\square$$

**Lemma A.3.**
$$\int_y U(X, y) dy = [f(X) + o(1)] V_u(X)^2.$$

*Proof.* We note that for $U(X, y) \neq 0$, we need $\{X, y\} \in \Psi_u^c$. We therefore have, $f(y) = f(X) + o(1)$.

$$
\begin{aligned}
\int_y U(X, y) dy &= \int [f(X) + o(1)] V(S_u(X) \cap S_u(Y)) dy \\
&= V_u(X)[f(X) + o(1)] \int \aleph(||X - y||/R_u(X)) dy \\
&= V_u(X)[f(X) + o(1)] R_u(X)^d \int \aleph(||y||/R_u(X)) d(y/R_u(X)) \\
&= V_u(X)[f(X) + o(1)] \frac{V_u(X)}{c_d} \int \aleph(||y||/R_u(X)) d(y/R_u(X)) \\
&= [f(X) + o(1)] \frac{V_u^2(X)}{c_d} \int \aleph(\delta) d(\delta).
\end{aligned}
$$

The integral $\int \aleph(\delta) d(\delta)$ can be shown to be equal to $c_d$ for all dimensions $d$.
We then have,

$$
\begin{aligned}
\int_y U(X, y) dy &= [f(X) + o(1)] V_u^2(X) \\
&= [f(X) + o(1)] \left( \frac{k}{M} \right)^2.
\end{aligned}
$$

$$\square$$

**Lemma A.4.** *Let* $\gamma_1(X)$, $\gamma_2(X)$ *be arbitrary continuous functions. Let* $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ *denote* $M + 2$ *i.i.d realizations of the density* $f$. *Then,*

$$
Cov[\gamma_1(\mathbf{X})\mathbf{e_u}(\mathbf{X}), \gamma_2(\mathbf{Y})\mathbf{e_u}(\mathbf{Y})] = \frac{Cov[\gamma_1(\mathbf{X})f(\mathbf{X}), \gamma_2(\mathbf{X})f(\mathbf{X})]}{M} + o(1/M).
$$

*Proof.*

$$
\begin{aligned}
Cov[\gamma_1(\mathbf{X})\mathbf{e_u}(\mathbf{X}), \gamma_2(\mathbf{Y})\mathbf{e_u}(\mathbf{Y})] &= \mathbb{E}\Big[\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})(\hat{\mathbf{f}}_\mathbf{u}(X) - \mathbb{E}[\hat{\mathbf{f}}_\mathbf{u}(X)])(\hat{\mathbf{f}}_\mathbf{u}(Y) - \mathbb{E}[\hat{\mathbf{f}}_\mathbf{u}(Y)])\Big] \\
&= \frac{1}{MV_u(X)V_u(Y)}\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})(U(\mathbf{X},\mathbf{Y}) - U(\mathbf{X})U(\mathbf{Y}))] \\
&= \frac{1}{MV_u^2(X)}\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})U(\mathbf{X},\mathbf{Y})] \\
&\quad - \frac{1}{MV_u^2(X)}\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})U(\mathbf{X})U(\mathbf{Y})] \\
&= I - II.
\end{aligned}
$$

$$
II = \frac{1}{M}\left(\mathbb{E}[\gamma_1(\mathbf{X})f(\mathbf{X})]\mathbb{E}[\gamma_2(\mathbf{Y})f(\mathbf{Y})]\right).
$$

$$
\begin{aligned}
I &= \frac{1}{MV_u^2(X)}\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})U(\mathbf{X},\mathbf{Y})] \\
&= \frac{1}{MV_u^2(X)}\int\int \gamma_1(x)\gamma_2(y)f(x)f(y)U(x,y)dxdy.
\end{aligned}
$$

Now for $U(x,y) \neq 0$, we need $\{x,y\} \in \Psi_u^c$. We therefore have, $\gamma_2(y)f(y) = \gamma_2(x)f(x) + o(1)$. We then have,

$$
\begin{aligned}
I &= \frac{1}{MV_u^2(X)}\int\int [\gamma_1(x)\gamma_2(x)f^2(x) + o(1)]U(x,y)dxdy \\
&= \frac{1}{MV_u^2(X)}\int [\gamma_1(x)\gamma_2(x)f^2(x) + o(1)]\left(\int U(x,y)dy\right)dx \\
&= \frac{1}{MV_u^2(X)}\int [\gamma_1(x)\gamma_2(x)f^2(x) + o(1)]\left((f(x) + o(1))V_u(x)^2\right)dx \\
&= \frac{1}{M}\int [\gamma_1(x)\gamma_2(x)f^2(x) + o(1)](f(x) + o(1))dx \\
&= \frac{1}{M}\left(\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})] + o(1)\right) \\
&= \frac{1}{M}\mathbb{E}[\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})] + o(1/M).
\end{aligned}
$$

$\square$

## A.6 Higher cross moments

**Disjoint balls** We have the following results concerning higher cross moments for disjoint balls:

**Lemma A.5.** *Let q,r be positive integers satisfying $q + r > 2$. For a fixed pair of points $\{X, Y\} \in \Psi_u{}^c$,*

$$Cov(\mathbf{e}_\mathbf{u}^q(X), \mathbf{e}_\mathbf{u}^r(Y)) \quad = \quad o(1/M).$$

*Proof.* For a fixed pair of points $\{X, Y\} \in \Psi_u{}^c$, the joint probability mass function of the functions $\mathbf{l_u}(X), \mathbf{l_u}(Y)$ is given by

$$Pr(\mathbf{l_u}(X) = l_x, \mathbf{l_u}(Y) = l_y) = 1_{l_x + l_y \le M} \binom{M}{l_x, l_y} (U(X))^{l_x} (U(Y))^{l_y} (1 - U(X) - U(Y))^{M - l_x - l_y}.$$

We also have from chernoff inequalities for binomial random variables that

$$Pr((1 - p)k < \mathbf{l_u}(X) < (1 + p)k) = 1 - e^{-p^2 k},$$
$$Pr((1 - p)k < \mathbf{l_u}(Y) < (1 + p)k) = 1 - e^{-p^2 k}.$$

Denote the high probability event $\chi$ by $(1 - p)k < \mathbf{l_u}(X), \mathbf{l_u}(Y) < (1 + p)k$. Define $\hat{\mathbf{l}}_\mathbf{u}(X)$, $\hat{\mathbf{l}}_\mathbf{u}(Y)$ to be binomial random variables with parameters $\{U(X), M - q\}$ and $\{U(Y), M - r\}$ respectively. The covariance between powers of density estimates is then given by

$$Cov(\hat{\mathbf{f}}_\mathbf{u}^q(X), \hat{\mathbf{f}}_\mathbf{u}^r(Y)) = \frac{1}{k^{q+r}} Cov(\mathbf{l}_\mathbf{u}^q(X), \mathbf{l}_\mathbf{u}^r(Y))$$

$$= \frac{1}{k^{q+r}} \sum l_x^q l_y^r Pr(\mathbf{l_u}(X) = l_x, \mathbf{l_u}(Y) = l_y) - \frac{1}{k^{q+r}} \sum l_x^q l_y^r Pr(\mathbf{l_u}(X) = l_x) Pr(\mathbf{l_u}(Y) = l_y)$$

$$= \sum_\chi \frac{l_x^q l_y^r}{k^{q+r}} \left[ Pr(\mathbf{l_u}(X) = l_x, \mathbf{l_u}(Y) = l_y) - Pr(\mathbf{l_u}(X) = l_x) Pr(\mathbf{l_u}(Y) = l_y) \right] + O(e^{-p^2 k})$$

$$= \sum_\chi \frac{f^q(X) f^r(Y) l_x^q l_y^r U^q(X) U^r(Y)}{k^{q+r}(l_x \times \ldots \times l_x - q + 1)(l_y \times \ldots \times l_y - r + 1)} \times$$

$$[(M \times \ldots \times M - (q + r - 1)) Pr(\hat{\mathbf{l}}_\mathbf{u}(X) = l_x, \hat{\mathbf{l}}_\mathbf{u}(Y) = l_y)$$

$$-(M \times \ldots \times M - q + 1)(M \times \ldots \times M - r + 1) Pr(\hat{\mathbf{l}}_\mathbf{u}(X) = l_x) Pr(\hat{\mathbf{l}}_\mathbf{u}(Y) = l_y)]$$

$$+ \quad o(1/M)$$

$$= \left( \frac{f^q(X) f^r(Y)}{M^{q+r}} + O\left( \frac{1}{k M^{q+r}} \right) \right) \times$$

$$\sum_\chi [(M \times \ldots \times M - (q + r - 1)) Pr(\hat{\mathbf{l}}_\mathbf{u}(X) = l_x, \hat{\mathbf{l}}_\mathbf{u}(Y) = l_y)$$

$$-(M \times \ldots \times M - (q - 1))(M \times \ldots \times M - (r - 1)) Pr(\hat{\mathbf{l}}_\mathbf{u}(X) = l_x) Pr(\hat{\mathbf{l}}_\mathbf{u}(Y) = l_y)]$$

$$+ \quad o(1/M)$$

$$
\begin{aligned}
&= \left( \frac{f^q(X)f^r(Y)}{M^{q+r}} + O\left( \frac{1}{kM^{q+r}} \right) \right) \times \\
&\quad [(M \times \ldots \times M - (q+r-1)) - (M \times \ldots \times M - (q-1))(M \times \ldots \times M - (r-1))] \\
&\quad + o(1/M) \\
&= \frac{-qr f^q(X)f^r(Y)}{M} + o\left( \frac{1}{M} \right).
\end{aligned}
$$

Then, the covariance between the powers of the error function is given by

$$
\begin{aligned}
Cov(e_{\mathbf{u}}^q(X), e_u^r(Y)) &= Cov((\hat{\mathbf{f}}_{\mathbf{u}}(X) - \mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(X)])^q, (\hat{\mathbf{f}}_{\mathbf{u}}(Y) - \mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(Y)])^r) \\
&= \sum_{a=1}^{q} \sum_{b=1}^{r} \binom{q}{a} \binom{r}{b} (-\mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(X)])^a (-\mathbb{E}[\hat{\mathbf{f}}_{\mathbf{u}}(Y)])^b Cov(\hat{\mathbf{f}}_{\mathbf{u}}^a(X), \hat{\mathbf{f}}_{\mathbf{u}}^b(Y)) \\
&= \sum_{a=1}^{q} \sum_{b=1}^{r} \binom{q}{a} \binom{r}{b} [(-f(X))^a(-f(Y))^b + o(1)] Cov(\hat{\mathbf{f}}_{\mathbf{u}}^a(X), \hat{\mathbf{f}}_{\mathbf{u}}^b(Y)) \\
&= -f^q(X)f^r(Y) \sum_{a=1}^{q} \sum_{b=1}^{r} \binom{q}{a} \binom{r}{b} \frac{(-1)^a a(-1)^b b}{M} + o\left( \frac{1}{M} \right) \\
&= 1_{\{q=1,r=1\}} \left( \frac{-f(X)f(Y)}{M} \right) + o(1/M) \\
&= o(1/M).
\end{aligned}
$$

where the last step follows from the condition that $q + r > 2$.

$\square$

**Intersecting balls**   For $\{X, Y\} \in \Psi_u{}^c$, we have the following bounds

**Lemma A.6.** *Let $\gamma_1(X)$, $\gamma_2(X)$ be arbitrary continuous functions. Let $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ denote $M + 2$ i.i.d realizations of the density $f$. Also let the indicator function $1_{\Delta_u}(X, Y)$ denote the event $\Delta_u : \{X, Y\} \in \Psi_u{}^c$. For $q,r$ positive integers satisfying $q + r > 1$,*

$$
\mathbb{E}[1_{\mathbf{\Delta_u}}(\mathbf{X}, \mathbf{Y}) \gamma_1(\mathbf{X}) \gamma_2(\mathbf{Y}) e_{\mathbf{u}}^q(\mathbf{X}) e_{\mathbf{u}}^r(\mathbf{Y})] = o\left( \frac{1}{M} \right),
$$

(55)

*Proof.* For $1_{\Delta_u}(X, Y) \neq 0$, we have $\{X, Y\} \in \Psi_u^c$. Then,

$$
\begin{aligned}
&\mathbb{E}[\mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y})\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e}_u^q(\mathbf{X})\mathbf{e}_u^r(\mathbf{Y})] \\
&= \mathbb{E}[\mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y})\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbb{E}_{\mathbf{X},\mathbf{Y}}[\mathbf{e}_u^q(X)\mathbf{e}_u^r(Y)]] \\
&\leq \mathbb{E}\left[\mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y})\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\sqrt{\mathbb{E}_{\mathbf{X}}[\mathbf{e}_u^{2q}(X)]\mathbb{E}_{\mathbf{Y}}[\mathbf{e}_u^{2r}(Y)]}\right] \\
&= \mathbb{E}\left[\mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y})\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})O\left(\frac{1}{k^{q+r/2}}\right)\right] \\
&= \int\left[O\left(\frac{1}{k^{q+r/2}}\right)(\gamma_1(x)\gamma_2(x) + o(1))\right]\left(\int \Delta_u(x, y)dy\right)dx \\
&= \int\left[O\left(\frac{1}{k^{q+r/2}}\right)(\gamma_1(x)\gamma_2(x) + o(1))\right]\left(2^d\frac{k}{M}\right)dx \\
&= o\left(\frac{1}{M}\right).
\end{aligned}
$$

where the bound is obtained using the Cauchy-Schwarz inequality and using Eq.51. $\qquad\square$

We can succinctly state the results derived in the last two lemmas in the form of the following lemma:

**Lemma A.7.** *Let $\gamma_1(X)$, $\gamma_2(X)$ be arbitrary continuous functions. Let $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ denote $M+2$ i.i.d realizations of the density $f$. If $q,r$ are positive integers satisfying $q+r > 2$*

$$
Cov[\gamma_1(\mathbf{X})\mathbf{e}_u^q(\mathbf{X}), \gamma_2(\mathbf{Y})\mathbf{e}_u^r(\mathbf{Y})] = o(1/M).
$$

*Proof.* The result for the case $q = 1$, $r = 1$ was established earlier in Lemma A.4.

$$
Cov[\gamma_1(\mathbf{X})\mathbf{e}_u^q(\mathbf{X}), \gamma_2(\mathbf{Y})\mathbf{e}_u^r(\mathbf{Y})] = I + D,
$$

where '$I$' stands for the contribution form the intersecting balls and '$D$' for the contribution from the dis-joint balls. $I$ and $D$ are given by

$$
\begin{aligned}
I &= \mathbb{E}[\mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y})Cov[\gamma_1(X)\mathbf{e}_u^q(X), \gamma_2(Y)\mathbf{e}_u^r(Y)]], \\
D &= \mathbb{E}[(\mathbf{1} - \mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y}))Cov[\gamma_1(X)\mathbf{e}_u^q(X), \gamma_2(Y)\mathbf{e}_u^r(Y)]].
\end{aligned}
$$

We have already established in the previous lemma that

$$
I = o\left(\frac{1}{M}\right).
$$

Now,

$$
\begin{aligned}
D &= \mathbb{E}[(\mathbf{1} - \mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y}))\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbb{E}_{\mathbf{X},\mathbf{Y}}[Cov(\mathbf{e}_u^q(X), \mathbf{e}_u^r(Y))]] \qquad (56) \\
&= \mathbb{E}[(\mathbf{1} - \mathbf{1}_{\Delta_u}(\mathbf{X}, \mathbf{Y}))\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})o(1/M)] \\
&= o\left(\frac{1}{M}\right).
\end{aligned}
$$

This concludes the proof. □

# B  $k$-NN density estimation

In this appendix, moment properties of the standard $k$-NN density estimate $\hat{\mathbf{f}}_k(X)$ are derived conditioned on $X_1, \ldots, X_N$. As the samples $X_1, \ldots, X_N, X_{N+1}, \ldots, X_T$, $T = M+N$ are i.i.d., these conditional moments are independent of the $N$ samples $\mathbf{X}_1, .., \mathbf{X}_N$.

## B.1  Preliminaries

Let $d(X, Y)$ denote the Euclidean distance between points $X$ and $Y$ and $\mathbf{d}_X^{(k)}$ denote the Euclidean distance between a point X and its $k$-th nearest neighbor amongst $\mathbf{X}_{N+1}, .., \mathbf{X}_{N+M}$. Let $c_d$ denote the unit ball volume in $d$ dimensions. The $k$-NN region is

$$\mathbf{S}_k(X) = \{Y : d(X, Y) \leq \mathbf{d}_X^{(k)}\}$$

and the volume of the $k$-NN region is

$$\mathbf{V}_k(X) = \int_{\mathbf{S}_k(X)} dZ.$$

The standard $k$-NN density estimator [30] is defined as

$$\hat{\mathbf{f}}_k(X) = \frac{k - 1}{M \mathbf{V}_k(X)}.$$

Define the coverage function as

$$\mathbf{P}(X) = \int_{\mathbf{S}_k(X)} f(Z) dZ.$$

Define spherical regions

$$S_r(X) = \{Y \in \mathbb{R}^d : d(X, Y) \leq r\}.$$

## B.2  Concentration inequality for coverage probability

It has been previously established that $\mathbf{P}(X)$ has a beta distribution with parameters $k$, $M - k + 1$. [31]. Consider a binomial random variable with parameters $M$ and $P$ with distribution function $Bi(.|M, P)$ and a beta random variable with parameters $k$ and $M - k + 1$ with distribution function $Be(.|k, M - k + 1)$. We have the following identity,

$$Be(P|k, M - k + 1) = 1 - Bi(k - 1|M, P). \tag{57}$$

The following Chernoff bounds for binomial random variables have also been established previously. When $k < MP$, $Bi(k|M, P) \leq exp[-(MP - k)^2/2PM]$, and when $k > MP$, $1 - Bi(k|M, P) \leq exp[-(MP - k)^2/2PM]$. We therefore have that for some $0 < p < 1/2$,

$$Pr((1-p)(k-1)/M < \mathbf{P}(X) < (p+1)(k-1)/M) = O(e^{-p^2k/2}). \tag{58}$$

Define

$$k_M = (k-1)/M.$$

Let $\natural(X)$ denote the event

$$\mathbf{P}(X) < (p_k + 1)k_M, \tag{59}$$

where $p_k = \sqrt{6}/(k^{\delta/2})$. Then, $1 - Pr(\natural(X)) = O(e^{-p_k^2 k/2}) = O(e^{-3k^{(1-\delta)}})$. Equivalently,

$$1 - Pr(\natural(X)) = O(\mathcal{C}(k)), \tag{60}$$

where $\mathcal{C}(k)$ is a function which satisfies the rate of decay condition $\mathcal{C}(k) = O(e^{-3k^{(1-\delta)}})$. Similarly, let $\natural_{-1}(X)$ denote the event

$$\mathbf{P}(X) > (1 - p_k)k_M, \tag{61}$$

Then

$$1 - Pr(\natural_{-1}(X)) = O(\mathcal{C}(k)), \tag{62}$$

Also let $\natural\natural(X) = \natural(X) \cap \natural_{-1}(X)$. Then

$$1 - Pr(\natural\natural(X)) = O(\mathcal{C}(k)), \tag{63}$$

Finally, we note that $\Gamma(x+a)/\Gamma(x) = x^a + o(x^a)$. Then for any $a < k$, $\mathbb{E}[\mathbf{P}^{-a}(X)]$ exists and is given by

$$\mathbb{E}[\mathbf{P}^{-a}(X)] = \frac{\Gamma(k-a)\Gamma(M+1)}{\Gamma(k)\Gamma(M+1-a)} = \Theta((k_M)^{-a}). \tag{64}$$

### B.2.1 Interior points

Let $\mathcal{S}'$ to be any arbitrary subset of $\mathcal{S}_I$ (2) satisfying the condition $Pr(\mathbf{Y} \notin \mathcal{S}') = o(1)$ where $\mathbf{Y}$ is random variable with density $f$. This implies that given the event $\natural(X)$, the $k$-NN neighborhoods $\mathbf{S}_k(X)$ of points $X \in \mathcal{S}'$ will lie completely inside the domain $\mathcal{S}$. Therefore the density $f$ has continuous partial derivatives of order $2\nu$ in the $k$-NN ball neighborhood $\mathbf{S}_k(X)$ for each $X \in \mathcal{S}'$ (assumption $(\mathcal{A}.2)$). We will now derive moments for the interior set of points $X \in \mathcal{S}'$. This excludes the set of points $X$ close to the boundary of the support whose $k$-NN neighborhoods $\mathbf{S}_k(X)$ intersect with the boundary of the support. We will deal with these points in Appendix B.

### B.2.2  Taylor series expansion of coverage probability

Let $X \in \mathcal{S}'$. Given the event $\natural(X)$, the coverage function $\mathbf{P}(X)$ can be represented in terms of the volume of the $k$-NN ball $\mathbf{V}_k(X)$ by expanding the density $f$ in a Taylor series about $X$ as follows. In particular, for some fixed $x \in \mathcal{S}'$, let

$$p(u) = \int_{S_u(x)} f(z) dz.$$

Using $(\mathcal{A}.2)$, we can write, by a Taylor series expansion of $f$ around $x$ using multi-index notation [39]

$$f(z) = \sum_{0 \le |\alpha| \le 2\nu} \frac{(z-x)^\alpha}{\alpha!}(\partial^\alpha f)(x) + o(||z - x||^{2\nu}) \tag{65}$$

Assuming $S_u(x) \subset \mathcal{S}$, we can then write

$$
\begin{aligned}
p(u) &= \int_{S_u(x)} f(z) dz \\
&= \int_{S_u(x)} \left( \sum_{|0 \le \alpha \le 2\nu|} \frac{(z-x)^\alpha}{\alpha!}(\partial^\alpha f)(x) \right) dz + o(u^{d+2\nu}) \\
&= f(x)c_d u^d + \sum_{i=1}^{\nu-1} c_i(x) c_d^{1+2i/d} u^{d+2i} + o(u^{d+2\nu}). 
\end{aligned}
\tag{66}
$$

where $c_i(x)$ are functionals of the derivatives of $f$. Now, denote $v(u) = \int_{S_u(x)} dz$ to be the volume of $S_u(x)$. Let $u^{inv}(v)$ be the inverse function of $v(u)$. Note that this inverse is well-defined since $v(u)$ is monotonic in $u$. Since $S_u(x) \subset \mathcal{S}$, $v(u) = c_d u^d$. This gives $u^{inv}(v) = (v/c_d)^{1/d}$. Define

$$P(v) = \int_{S_{u^{inv}(v)}(x)} f(z) dz.$$

Using (66),

$$P(v) = f(X)v + \sum_{i=1}^{\nu-1} c_i(X) v^{1+2i/d} + o(v^{1+2\nu/d}). \tag{67}$$

Now denote $V(p) = P^{inv}(p)$ to be the inverse of $P(.)$. Note that this inverse is well-defined since $P(v)$ is monotonic in $v$. Dividing (67) by $vP(v)$ on both sides, we get

$$\frac{1}{v} = \frac{f(X)}{P(v)} + \sum_{i=1}^{\nu-1} \frac{c_i(X)}{P(v)} v^{2i/d} + o(v^{2\nu/d} P^{-1}(v)) \tag{68}$$

By repeatedly substituting the LHS of (68) in the RHS of (68), we can obtain (69):

$$\frac{1}{V(p)} = \frac{f(X)}{p} + \sum_{i=1}^{\nu-1} \frac{h_i(X)}{p^{1-2i/d}} + o(p^{2\nu/d-1}), \tag{69}$$

From our derivation of (69) using (67), it is clear that $h_i(X)$ are of the form

$$h_i(X) = \sum_{\{a_i\}=A; A \in \mathcal{A}} \frac{\prod_{i=1}^{\nu-1} c_i^{a_i}}{f^{a_0}(X)}$$

where $A$ is a $\nu$-tuple of positive real numbers $a_0, .., a_{\nu-1}$ and the cardinality of $\mathcal{A}$ is finite. By assumptions ($\mathcal{A}$.1) and ($\mathcal{A}$.2), this implies that the constants $h_i(X)$ are *bounded*. Also, we note that $h(X) = h_1(X) = c(X)f^{-2/d}(X)$ [15], where $c(X) := c_1(X) = \Gamma^{(2/d)}(\frac{d+2}{2})tr[\nabla^2(f(X))]$. This then implies that under the event $\natural(X)$

$$\frac{1}{\mathbf{V}_k(X)} = \frac{f(X)}{\mathbf{P}(X)} + \sum_{t \in \mathcal{T}} \frac{h_t(X)}{\mathbf{P}^{1-t}(X)} + \mathbf{h_r}(X), \tag{70}$$

where $\mathcal{T} = \{2/d, 4/d, 6/d.., 2\nu/d\}$ and $\mathbf{h_r}(X) = o(\mathbf{P}^{2\nu/d-1}(X))$. Now, by ($\mathcal{A}$.2), we have $(k/M)^{2\nu/d} = o(1/M)$. This implies that $2\nu/d > 1$. Under the event $\natural(X)$, we have $\mathbf{P}(X) \leq (p_k + 1)k/M$, which, in conjunction with the condition $2\nu/d > 1$ implies that

$$\mathbf{h_r}(X) = o(\mathbf{P}^{2\nu/d-1}(X)) = o((k/M)^{2\nu/d-1}) = o(1/k_M M). \tag{71}$$

On the other hand, under the event, $\natural^c(X)$, $(p_k + 1)k/M \leq \mathbf{P}(X) \leq 1$, which gives

$$\mathbf{h_r}(X) = O(1). \tag{72}$$

### B.2.3 Approximation to the $k$-NN density estimator

Define the *coverage* density estimate to be,

$$\hat{\mathbf{f}}_c(X) = f(X)\frac{k-1}{M}\frac{1}{\mathbf{P}(X)}.$$

The estimate $\hat{\mathbf{f}}_c(X)$ is clearly not implementable. Note also that the two estimates - $\hat{\mathbf{f}}_c(X)$ and $\hat{\mathbf{f}}_k(X)$ - are identical in the case of the uniform density.

$$\frac{1}{\mathbf{V}_k(X)} = \frac{f(X)}{\mathbf{P}(X)} + \frac{h(X)}{\mathbf{P}^{1-2/d}(X)} + \mathbf{h_s}(X), \tag{73}$$

where $\mathbf{h_s}(X) = o(1/\mathbf{P}^{1-2/d}(X))$. This gives,

$$\hat{\mathbf{f}}_k(X) = \hat{\mathbf{f}}_c(X) + \left(\frac{k-1}{M}\right)\frac{h(X)}{\mathbf{P}^{1-2/d}(X)} + \frac{k-1}{M}\mathbf{h_s}(X). \tag{74}$$

whenever $\natural(X)$ is true.

56

### B.2.4 Bounds on $k$-NN density estimates

Let $X$ be a Lebesgue point of $f$, i.e., an $X$ for which

$$\lim_{r \to 0} \frac{\int_{S_r(X)} f(y)dy}{\int_{S_r(x)} dy} = f(X).$$

Because $f$ is an density, we know that almost all $X \in \mathcal{S}$ satisfy the above property. Now, fix $\epsilon \in (0, 1)$ and find $\delta > 0$ such that

$$\sup_{0 < r \leq \delta} \frac{\int_{S_r(X)} f(y)dy}{\int_{S_r(x)} dy} - f(X) \leq \epsilon f(X).$$

This in turn implies that, for $\mathbf{P}(X) \leq P(\delta)$,

$$\frac{\mathbf{P}(X)}{(1 + \epsilon)f(X)} \leq \mathbf{V}_k(X) \leq \frac{\mathbf{P}(X)}{(1 - \epsilon)f(X)} \tag{75}$$

and in turn implies

$$(1 - \epsilon)\hat{\mathbf{f}}_c(X) \leq \quad \hat{\mathbf{f}}_k(X) \quad \leq (1 + \epsilon)\hat{\mathbf{f}}_c(X). \tag{76}$$

Also, because $\delta > 0$ is fixed, we note that the event $\mathbf{P}(X) \leq P(\delta)$ is a subset of $\natural(X)$ and therefore (75) holds under $\natural(X)$.

Under the event $\natural^c(X)$, we can bound $\mathbf{V}_k(X)$ from above by $c_d \mathcal{D}^d$. Also, since $\mathbf{V}_k(X)$ is monotone in $\mathbf{P}(X)$, under the event $\natural^c(X)$, we can bound $\mathbf{V}_k(X)$ from below by $(1 + p_k)(k-1)/M(1 - \epsilon)f(X)$ and therefore by $(k-1)/M(1 - \epsilon)f(X)$. Written explicitly,

$$\frac{(k-1)}{M(1 - \epsilon)f(X)} \leq \mathbf{V}_k(X) \leq c_d \mathcal{D}^d \tag{77}$$

and in turn implies

$$(k-1)/(Mc_d \mathcal{D}^d) \leq \quad \hat{\mathbf{f}}_k(X) \quad \leq (1 - \epsilon)f(X). \tag{78}$$

Finally, note that $k_M/\mathbf{P}(X)$ is bounded above by $O(1)$ under the event $\natural(X)$. This implies that for any $a < k$,

$$\mathbb{E}[\natural^c(X)]k_M^a \mathbf{P}^{-a}(X) \leq O(1)Pr(\natural^c(X)) = O(\mathcal{C}(k)). \tag{79}$$

## B.3   Bias of the $k$-NN density estimates

Let $X \in \mathcal{S}'$. We can analyze the bias of $k$-NN density estimates as follows by using (74)

$$
\begin{aligned}
\mathbb{E}[1_{\natural(X)}\hat{\mathbf{f}}_k(X)] \;=\;& \mathbb{E}[1_{\natural(X)}\hat{\mathbf{f}}_c(X)] + \mathbb{E}\left[1_{\natural(X)}\left(\frac{k-1}{M}\right)\frac{h(X)}{\mathbf{P}^{1-2/d}(X)}\right] + \mathbb{E}\left[1_{\natural(X)}\frac{k-1}{M}\mathbf{h_s}(X)\right] \\
=\;& \mathbb{E}[1_{\natural(X)}\hat{\mathbf{f}}_c(X)] + \mathbb{E}\left[1_{\natural(X)}\left(\frac{k-1}{M}\right)\frac{h(X)}{\mathbf{P}^{1-2/d}(X)}\right] + o\left(\mathbb{E}\left[1_{\natural(X)}\frac{k-1}{M}\mathbf{P}^{2/d-1}(X)\right]\right) \\
=\;& \mathbb{E}[\hat{\mathbf{f}}_c(X)] + \mathbb{E}\left[\left(\frac{k-1}{M}\right)\frac{h(X)}{\mathbf{P}^{1-2/d}(X)}\right] + o\left(\frac{k}{M}\right)^{2/d} + O(\mathcal{C}(k)) \\
=\;& f(X) + h(X)\left(\frac{k}{M}\right)^{2/d} + o\left(\frac{k}{M}\right)^{2/d},
\end{aligned} \tag{80}
$$

where we used the fact that under the event $\natural^c(X)$, $((k-1)/M)\mathbf{P}^{1-t}(X) = O(1)$ for any $t >= 0$, which in turn gives $\mathbb{E}[1_{\natural^c(X)}((k-1)/M)\mathbf{P}^{1-t}(X)] = O(Pr(\natural^c(X))) = O(\mathcal{C}(k))$. This implies that

$$
\begin{aligned}
\mathbb{E}[\hat{\mathbf{f}}_k(X)] - f(X) \;=\;& \mathbb{E}[1_{\natural(X)}\hat{\mathbf{f}}_k(X)] + \mathbb{E}[1_{\natural^c(X)}\hat{\mathbf{f}}_k(X)] - f(X) \\
=\;& h(X)\left(\frac{k}{M}\right)^{2/d} + o\left(\frac{k}{M}\right)^{2/d} + O(\mathcal{C}(k)) + \mathbb{E}[1_{\natural^c(X)}\hat{\mathbf{f}}_k(X)] \\
=\;& h(X)\left(\frac{k}{M}\right)^{2/d} + o\left(\frac{k}{M}\right)^{2/d} + O(\mathcal{C}(k)),
\end{aligned} \tag{81}
$$

where the last step follows because , by (78), $1_{\natural^c(X)}\hat{\mathbf{f}}_k(X) = O(1)$. This expression is true for $k >= 3$ by (64).

Next, assuming that (8) holds, we evaluate $\mathbb{E}[g(\hat{\mathbf{f}}_k(X), X)]$ in an identical fashion to the derivation of (81).

$$
\begin{aligned}
\mathbb{E}[1_{\natural(X)}g(\hat{\mathbf{f}}_k(X), X)] &= \mathbb{E}\left[1_\natural(X)g\left(\hat{\mathbf{f}}_c(X) + k_M h(X)(\mathbf{P}(X))^{2/d-1} + k_M \mathbf{h_s}(X), X\right)\right] \\
&= \mathbb{E}\left[1_{\natural(X)}g\left(\hat{\mathbf{f}}_c(X) + k_M h(X)(\mathbf{P}(X))^{2/d-1} + k_M o((\mathbf{P}(X))^{2/d-1}), X\right)\right] \\
&= \mathbb{E}\left[g\left(\hat{\mathbf{f}}_c(X) + k_M h(X)(\mathbf{P}(X))^{2/d-1} + k_M o((\mathbf{P}(X))^{2/d-1}), X\right)\right] + O(\mathcal{C}(k)) \\
&= \mathbb{E}\left[g(\hat{\mathbf{f}}_c(X), X) + g'(\hat{\mathbf{f}}_c(X), X)k_M h(X)(\mathbf{P}(X))^{2/d-1} + o(k_M \mathbf{P}(X))^{2/d-1})\right] + O(\mathcal{C}(k)) \\
&= g(f(X), X)g_1(k, M) + g_2(k, M) + g'(f(X), X)h(X)(k/M)^{2/d} + o((k/M)^{2/d}) + O(\mathcal{C}(k)).
\end{aligned}
$$

This gives,

$$
\begin{aligned}
\mathbb{E}[g(\hat{\mathbf{f}}_k(X), X)] &= \mathbb{E}[1_{\natural(X)}g(\hat{\mathbf{f}}_k(X), X)] + \mathbb{E}[1_{\natural^c(X)}g(\hat{\mathbf{f}}_k(X), X)] \\
&= g(f(X), X)g_1(k, M) + g_2(k, M) + g'(f(X), X)h(X)(k/M)^{2/d} + o((k/M)^{2/d}) + O(\mathcal{C}(k)).
\end{aligned} \tag{82}
$$

58

## B.4 Moments of error function

Let $\gamma_1(X)$, $\gamma_2(X)$ be arbitrary continuous functions satisfying the condition: $\sup_X[\gamma_i(X)]$ is finite, $i = 1, 2$. Also let $\gamma(X) = \gamma_1(X)$. Let $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ denote $M + 2$ i.i.d realizations of the density $f$. Let $q, r$ be arbitrary positive integers less than $k$. Define the error function

$$\mathbf{e}_k(X) = \hat{\mathbf{f}}_k(X) - \mathbb{E}[\hat{\mathbf{f}}_k(X) \mid X].$$

Then,

**Lemma B.1.**

$$\mathbb{E}\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^q(\mathbf{X})\right] = O(k^{-q\delta/2}) + o(1/M) + O(\mathcal{C}(k)). \tag{83}$$

**Lemma B.2.**

$$\begin{aligned}
Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k^q(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k^r(\mathbf{Y})\right] &= O\left(\frac{1}{k^{((q+r)\delta/2-1)}M}\right) + O(k_M^{2/d}/M) \\
&\quad + O(1/M^2) + O(\mathcal{C}(k)).
\end{aligned} \tag{84}$$

Define the operator $\mathcal{M}(\mathbf{Z}) = \mathbf{Z} - \mathbb{E}[\mathbf{Z}]$. Let $\beta$ be any positive real number and define

$$\mathbf{E}_\beta(X) = k_M^\beta(\mathcal{M}(\mathbf{P}^{-\beta}(X))). \tag{85}$$

Define the terms

$$\mathbf{e}_c(X) = \hat{\mathbf{f}}_c(X) - \mathbb{E}[\hat{\mathbf{f}}_c(X) \mid X], \tag{86}$$

$$\mathbf{e}_t(X) = \mathcal{M}\left(\sum_{t\in\mathcal{T}}\frac{k_M h_t(X)}{\mathbf{P}^{1-t}(X))}\right), \tag{87}$$

$$\mathbf{e}_r(X) = \mathcal{M}(k_M \mathbf{h_r}(X)). \tag{88}$$

Note that

$$\mathbf{e}_c(X) = f(X)\mathbf{E}_1(X) \tag{89}$$

and

$$\mathbf{e}_t(X) = (\sum_{t\in\mathcal{T}}k_M^t h_t(X)(\mathbf{E}_{1-t}(X))). \tag{90}$$

Define the event $\{X \in \mathcal{S}'\} \cap \{\natural(X)\}$ by $\dagger(X)$. Note that under the event $\dagger(X)$, $\mathbf{e}_k(X) = \mathbf{e}_c(X) + \mathbf{e}_t(X) + \mathbf{e}_r(X) =: \mathbf{e}_o(X)$. Also, under the event $\natural(X)$, $\mathbf{P}(X) \leq (1 + p_k)k_M$, which implies that under the event $\natural(X)$, the following hold

$$\mathbf{E}_\beta(X) = O(1), \mathbf{e}_c(X) = O(1), \mathbf{e}_t(X) = O(1), \mathbf{e}_r(X) = O(1), \mathbf{e}_o(X) = O(1). \tag{91}$$

Furthermore, by (78), under the event $\natural(X)$,

$$\mathbf{e}_k(X) = O(1). \tag{92}$$

59

*Proof.* of Lemma B.1. Since $\mathbf{P}(X)$ is a beta random variable, the probability density function of $\mathbf{P}(X)$ is given by

$$f(p_X) = \frac{M!}{(k-1)!(M-k)!}p_X^{k-1}(1-p_X)^{M-k}.$$

By (64), $\mathbb{E}[\mathbf{P}^{-\beta}(X)] = \Theta((k/M)^{-\beta})$ if $\beta < k$. We will first show that $\mathbb{E}[\mathbf{E}_\beta^q(X)] = O(1)$ if $q\beta < k$. This in turn implies that, by (89) and (90), $\mathbb{E}[\mathbf{e}_c^q(X)] = O(1)$ and $\mathbb{E}[\mathbf{e}_t^q(X)] = O(1)$ for any $q < k$.

$$\begin{aligned}
\mathbb{E}[\mathbf{E}_\beta^q(X)] &= \mathbb{E}\left[k_M^{q\beta}(\mathbf{P}^{-\beta}(X) - \mathbb{E}[\mathbf{P}^{-\beta}(X)])^q\right] \\
&= k_M^{q\beta}\sum_{i=1}^q \binom{q}{i}(-1)^{q-i}\mathbb{E}[\mathbf{P}^{-i\beta}(X)]\mathbb{E}[\mathbf{P}^{-(q-i)\beta}(X)] \\
&= k_M^{q\beta}\sum_{i=1}^q \binom{q}{i}(-1)^{q-i}\Theta((k/M)^{-i\beta})\Theta((k/M)^{-(q-i)\beta}) \\
&= \sum_{i=1}^q \binom{q}{i}(-1)^{q-i}\Theta(1) = O(1).
\end{aligned} \tag{93}$$

By (63) and (93),

$$\mathbb{E}[1_{\natural\natural^c(X)}\mathbf{E}_\beta^q(X)] = O(\mathcal{C}(k)).$$

By the definition of $\natural\natural(X)$,

$$1_{\natural\natural(X)}\mathbf{E}_\beta^q(X) = O\left(k^{-(\delta q/2)}\right), \tag{94}$$

and therefore

$$\mathbb{E}[1_{\natural\natural(X)}\mathbf{E}_\beta^q(X)] = O\left(k^{-(\delta q/2)}\right).$$

This gives,

$$\mathbb{E}[\mathbf{E}_\beta^q(X)] = O(k^{-\delta q/2}) + O(\mathcal{C}(k)). \tag{95}$$

From this analysis on $\mathbf{E}_\beta(X)$, it trivially follows from (89) that

$$\mathbb{E}[\mathbf{e}_c^l(X)] = O(k^{-\delta l/2}) + O(\mathcal{C}(k)). \tag{96}$$

Also observe that by (71) and (72),

$$\mathbb{E}[\mathbf{e}_r^l(X)] = \mathbb{E}[1_{\natural(X)}\mathbf{e}_r^l(X)] + \mathbb{E}[1_{\natural^c(X)}\mathbf{e}_r^l(X)] = o(1/M^l) + O(\mathcal{C}(k)). \tag{97}$$

We will now bound $\mathbf{e}_t^l(X)$. Let $L = \sum_{t\in\mathcal{T}} l_t t$. Now, using (90), $\mathbf{e}_t^l(X)$ can be expressed as a *sum* of terms of the form $(k/M)^L \binom{l}{l_1,..,l_t}\prod_{t\in\mathcal{T}}(h_t^l(X)\mathbf{E}_t^{l_t}(X))$ where $\sum_t l_t = l$. Now, we can bound each of these summands using (94) as follows:

$$\begin{aligned}
(k/M)^l\mathbb{E}[\prod_{t\in\mathcal{T}}\mathbf{E}_t^{l_t}(X)] &= (k/M)^L\mathbb{E}[1_{\natural\natural(X)}\prod_{t\in\mathcal{T}}\mathbf{E}_t^{l_t}(X)] + (k/M)^L\mathbb{E}[1_{\natural\natural^c(X)}\prod_{t\in\mathcal{T}}\mathbf{E}_t^{l_t}(X)] \\
&= (k/M)^L\prod_{t\in\mathcal{T}}O(k^{-l_t\delta/2}) + O(\mathcal{C}(k)) \\
&= (k/M)^L O(k^{-l\delta/2}) + O(\mathcal{C}(k)) \\
&= o(k^{-l\delta/2}) + O(\mathcal{C}(k)).
\end{aligned} \tag{98}$$

60

This implies that
$$\mathbb{E}[\mathbf{e}_t^l(X)] = o(k^{-l\delta/2}) + O(\mathcal{C}(k)). \tag{99}$$

Note that $\mathbf{e}_o^q(X)$ will contain terms of the form $(\mathbf{e}_c(X) + \mathbf{e}_t(X))^l(\mathbf{e}_r(X))^{q-l}$. If $l < q$, the expectation of this term can be bounded as follows

$$
\begin{aligned}
&|\mathbb{E}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))^l(\mathbf{e}_r(X))^{q-l}]| \\
&\leq \sqrt{\mathbb{E}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))^{2l}]\mathbb{E}[(\mathbf{e}_r(X))^{2(q-l)}]} \\
&= \sqrt{O(1)^{2l}(o(1/M))^{2(q-l)}} \\
&= O(1) \times (o(1/M))^{q-l} = o(1/M).
\end{aligned} \tag{100}
$$

Let us concentrate on the case $l = q$. In this case, $\mathbf{e}_k^q(X)$ will contain terms of the form $(\mathbf{e}_c(X))^m(\mathbf{e}_t(X))^{q-m}$. For $m < q$,

$$
\begin{aligned}
&|\mathbb{E}[(\mathbf{e}_c(X))^m(\mathbf{e}_t(X))^{q-m}]| \\
&\leq \sqrt{\mathbb{E}[(\mathbf{e}_c(X))^{2l}]\mathbb{E}[(\mathbf{e}_t(X))^{2(q-l)}]} \\
&= \left(O(k^{-m\delta/2}) \times o(k^{-(q-m)\delta/2})\right) + \mathcal{C}(k) = o(k^{-q\delta/2}) + O(\mathcal{C}(k)).
\end{aligned} \tag{101}
$$

This therefore implies that, by (96), (97), (99), (100) and (101),

$$
\begin{aligned}
\mathbb{E}[\mathbf{e}_o^q(X)] &= \mathbb{E}[\mathbf{e}_c^q(X)] + o(k^{-q\delta/2}) + \mathcal{C}(k) \\
&= O(k^{-q\delta/2}) + o(k^{-q\delta/2}) + o(1/M) + \mathcal{C}(k) \\
&= O(k^{-q\delta/2}) + o(1/M) + \mathcal{C}(k).
\end{aligned} \tag{102}
$$

This finally implies that

$$
\begin{aligned}
\mathbb{E}\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^q(\mathbf{X})\right] &= \mathbb{E}\left[1_{\dagger(\mathbf{x})}\gamma(\mathbf{X})\mathbf{e}_k^q(\mathbf{X})\right] + O(\mathcal{C}(k)) \quad (by(92)) \\
&= \mathbb{E}\left[1_{\dagger(\mathbf{x})}\gamma(\mathbf{X})\mathbf{e}_o^q(\mathbf{X})\right] + O(\mathcal{C}(k)) \\
&= \mathbb{E}\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_o^q(\mathbf{X})\right] + O(\mathcal{C}(k)) \quad (by(91)) \\
&= O(k^{-q\delta/2}) + o(1/M) + O(\mathcal{C}(k)).
\end{aligned} \tag{103}
$$

This concludes the proof.

$\square$

Before proving Lemma B.2, we seek to answer the following question: for which set of pair of points $\{X, Y\}$ are the $k$-NN balls disjoint?

### B.4.1 Intersecting and disjoint balls

Define $\Psi_\epsilon := \{X, Y\} \in \mathcal{S}' : ||X - Y|| \geq R_\epsilon(X) + R_\epsilon(Y)$ where $R_\epsilon(X)$ and $R_\epsilon(Y)$ are the ball radii of the spherical regions $S_u(X)$ and $S_u(Y)$, such that $\int_{S_u(X)} f(z)dz = \int_{S_u(Y)} f(z)dz = (1 + p_k)k_M$. We will now show that for $\{X, Y\} \in \Psi_\epsilon$, the $k$-NN balls will be disjoint with exponentially high probability. Let $\mathbf{d_X^{(k)}}$ and $\mathbf{d_Y^{(k)}}$ denote the $k$-NN distances from $X$ and $Y$ and let $\mathbf{\Upsilon}$ denote the event that the $k$-NN balls intersect. For $\{X, Y\} \in \Psi_\epsilon$,

$$
\begin{aligned}
Pr(\mathbf{\Upsilon}) &= Pr(\mathbf{d_X^{(k)}} + \mathbf{d_Y^{(k)}} \geq ||X - Y||) \\
&\leq Pr(\mathbf{d_X^{(k)}} + \mathbf{d_Y^{(k)}} \geq R_\epsilon(X) + R_\epsilon(Y)). \\
&\leq Pr(\mathbf{d_X^{(k)}} \geq R_\epsilon(X)) + Pr(\mathbf{d_Y^{(k)}} \geq R_\epsilon(Y)) \\
&= Pr(\mathbf{P}(X) \geq (p_k + 1)((k-1)/M)) \\
&\quad + Pr(\mathbf{P}(Y) \geq (p_k + 1)((k-1)/M)) \\
&= 2\mathcal{C}(k),
\end{aligned}
$$

where the last inequality follows from the concentration inequality (58). We conclude that for $\{X, Y\} \in \Psi_\epsilon$, the probability of intersection of $k$-NN balls centered at $X$ and $Y$ decays exponentially in $p_k^2 k$. Stated in a different way, we have shown that for a given pair of points $\{X, Y\}$, if the $\epsilon$ balls around these points are disjoint, then the $k$-NN balls will be disjoint with exponentially high probability. Let $\Delta_\epsilon(X, Y)$ denote the event $\{X, Y\} \in \Psi_\epsilon^c$. From the definition of the region $\Psi_\epsilon$, we have $Pr(\{\mathbf{X}, \mathbf{Y}\} \in \Psi_\epsilon^c) = O(k/M)$.

Let $\{X, Y\} \in \Psi_\epsilon$ and let $q, r$ be non-negative integers satisfying $q + r > 1$. The event that the $k$-NN balls intersect is given by $\mathbf{\Upsilon} := \{\mathbf{d_X^{(k)}} + \mathbf{d_Y^{(k)}} > ||X - Y||\}$. The joint probability distribution of $\mathbf{P}(X)$ and $\mathbf{P}(Y)$ when the $k$-NN balls do not intersect $=: \mathbf{\Upsilon^c}$ is given by

$$
f_{\mathbf{\Upsilon^c}}(p_X, p_Y) = M! \frac{(p_X p_Y)^{k-1}}{(k-1)!^2} \frac{(1 - p_X - p_Y)^{M-2k}}{(M - 2k)!}.
$$

Define

$$
i(p_X, p_Y) = \frac{\Gamma(t)\Gamma(u)\Gamma(v)}{\Gamma(t + u + v)} p_X^{t-1} p_Y^{u-1}(1 - p_X - p_Y)^{v-1},
$$

and note that

$$
\int_{p_X=0}^{1} \int_{p_Y=0}^{1} 1_{\{p_X + p_Y \leq 1\}} i(p_X, p_Y) dp_X dp_Y = 1.
$$

Figure 22 shows the distribution of the $M$ samples when the $k$-NN balls are disjoint. Now note that $i(p_X, p_Y)$ corresponds to the density function $f_{\mathbf{\Upsilon^c}}(p_X, p_Y)$ for the choices $t = k$, $u = k$ and $v = M - 2k + 1$. Furthermore, for $\{X, Y\} \in \Psi_\epsilon$, the set $\mathcal{Q} := \{p_X, p_Y\} : p_X, p_Y \leq (1 + p_k)(k-1)/M$ is a subset of the region $\mathcal{T} := \{p_X, p_Y\} : 0 \leq p_X, p_Y \leq 1; p_X + p_Y \leq 1$. Note that $\mathbb{E}[1_\mathcal{Q}] = 1 - \mathcal{C}(k)$. This implies that expectations over the region $\mathcal{R} := \{p_X, p_Y\} : 0 \leq p_X, p_Y \leq 1;$ should be of the same order as the expectations over $\mathcal{T}$ with differences of order $\mathcal{C}(k)$. In particular, for $t, u < k$,

$$
\mathbb{E}[\mathbf{P}^{-t}(X)\mathbf{P}^{-u}(Y)] = \mathbb{E}[1_\mathcal{T}\mathbf{P}^{-t}(X)\mathbf{P}^{-u}(Y)] + \mathcal{C}(k).
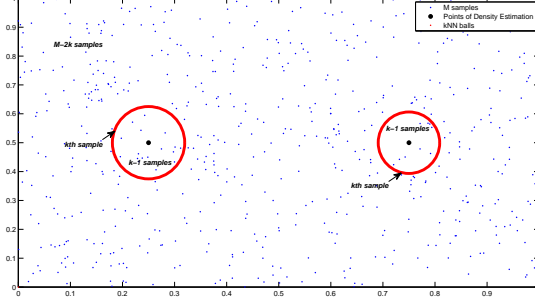$$

Figure 22: Distribution of samples when $k$-NN balls are disjoint.

From the joint distribution representation, it follows that

$$\frac{\mathbb{E}[1_{\mathcal{J}}\mathbf{P}^{-t}(X)\mathbf{P}^{-u}(Y)]}{\mathbb{E}[\mathbf{P}^{-t}(X)]\mathbb{E}[\mathbf{P}^{-u}(Y)]} = \frac{\Gamma(M-t)\Gamma(M-u)}{\Gamma(M-t-u)\Gamma(M)} = -\frac{tu}{M} + O(1/M^2). \tag{104}$$

Now observe that

$$
\begin{aligned}
&(k_M)^{t+u}Cov(\mathbf{P}^{-t}(X),\mathbf{P}^{-u}(Y)) \\
&= (k_M)^{t+u}[\mathbb{E}[\mathbf{P}^{-t}(X)\mathbf{P}^{-u}(Y)] - \mathbb{E}[\mathbf{P}^{-t}(X)]\mathbb{E}[\mathbf{P}^{-u}(Y)]] \\
&= (k_M)^{t+u}\mathbb{E}[\mathbf{P}^{-t}(X)]\mathbb{E}[\mathbf{P}^{-u}(Y)] \left[ \frac{\mathbb{E}[\mathbf{P}^{-t}(X)\mathbf{P}^{-u}(Y)]}{\mathbb{E}[\mathbf{P}^{-t}(X)]\mathbb{E}[\mathbf{P}^{-u}(Y)]} - 1 \right] \\
&= (k_M)^{t+u}\Theta(k_M^{-t})\Theta(k_M^{-u}) \left[ 1 - \frac{tu}{M} + o(1/M^2) - 1 \right] \qquad \text{(by (64) and (104))} \\
&= -\left(\frac{tu}{M}\right) + O(1/M^2). \tag{105}
\end{aligned}
$$

Then, the covariance between the powers of the error function $\mathbf{E}_\beta$, for $qt, ru < k$ is given by

$$
\begin{aligned}
Cov(\mathbf{E}_t^q(X),\mathbf{E}_u^r(Y)) &= k_M^{(tq+ur)}Cov\left(\left[\mathbf{P}^{-t}(X) - \mathbb{E}\left[\mathbf{P}^{-t}(X)\right]\right]^q, \left[\mathbf{P}^{-u}(Y) - \mathbb{E}\left[\mathbf{P}^{-u}(Y)\right]\right]^r\right) \\
&= \sum_{a=1}^{q}\sum_{b=1}^{r}\binom{q}{a}\binom{r}{b}[(-1)^{a+b} + o(1)]k_M^{(ta+ub)}Cov(\mathbf{P}^{-ta}(X),\mathbf{P}^{-ub}(Y)) \\
&= -tu\sum_{a=1}^{q}\sum_{b=1}^{r}\binom{q}{a}\binom{r}{b}\frac{(-1)^a a(-1)^b b}{M} + O\left(\frac{1}{M^2}\right) \\
&= 1_{\{q=1,r=1\}}\left(\frac{-tu}{M}\right) + O(1/M^2). \tag{106}
\end{aligned}
$$

*Proof.* of Lemma B.2. Let $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ denote $M+2$ i.i.d realizations of the density $f$.

63

Then, identical to the derivation of (103) in the proof of Lemma B.1,

$$Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k^q(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k^r(\mathbf{Y})\right]$$
$$= Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_o^q(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_o^r(\mathbf{Y})\right] + O(\mathcal{C}(k)).$$

Using the exact same arguments as in proof of Lemma A.1, it can be shown that the contribution of terms $\mathbf{e}_r(\mathbf{X}), \mathbf{e}_r(\mathbf{Y})$ to the R.H.S. of the above equation is $o(1/M)$. Define $\sharp(\mathbf{X},\mathbf{Y}) := \gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})Cov_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X)+\mathbf{e}_t(X))^q, (\mathbf{e}_c(Y)+\mathbf{e}_t(Y))^r]$. Thus,

$$Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k^q(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k^r(\mathbf{Y})\right]$$
$$= \mathbb{E}[1_{\{\mathbf{X},\mathbf{Y}\in\mathcal{S}'\}}\sharp(\mathbf{X},\mathbf{Y})] + O(\mathcal{C}(k))$$
$$= \mathbb{E}[1_{\mathbf{\Delta}_\epsilon^c(\mathbf{X},\mathbf{Y})}\sharp(\mathbf{X},\mathbf{Y})] + \mathbb{E}[1_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\sharp(\mathbf{X},\mathbf{Y})] + O(\mathcal{C}(k))$$
$$= I + II + O(\mathcal{C}(k)).$$

**For $\{X,Y\} \in \Psi_\epsilon^c$** The covariance term $Cov_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X)+\mathbf{e}_t(X))^q, (\mathbf{e}_c(Y)+\mathbf{e}_t(Y))^r]$ can be shown to be $O(k^{-(q+r)\delta/2})$ for $q, r < k$ by using Cauchy-Schwarz and (100), (101) as follows.

$$|Cov[(\mathbf{e}_c(X)+\mathbf{e}_t(X))^q, (\mathbf{e}_c(Y)+\mathbf{e}_t(Y))^r]| \leq \sqrt{\mathbb{V}[(\mathbf{e}_c(X)+\mathbf{e}_t(X))^q]\mathbb{V}[(\mathbf{e}_c(Y)+\mathbf{e}_t(Y))^r]}$$
$$\leq \sqrt{\mathbb{E}[(\mathbf{e}_c(X)+\mathbf{e}_t(X))^{2q}]\mathbb{E}[(\mathbf{e}_c(Y)+\mathbf{e}_t(Y))^{2r}]}$$
$$= \sqrt{O(k^{-(2q)\delta/2})O(k^{-(2r)\delta/2})}$$
$$= O(k^{-(q+r)\delta/2}). \tag{107}$$

This implies that

$$II = \mathbb{E}[1_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\sharp(\mathbf{X},\mathbf{Y})] = \mathbb{E}\left[1_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}O(k^{-(q+r)\delta/2})\right] = O\left(\frac{1}{k^{((q+r)\delta/2-1)}M}\right),$$

where the last but one step follows since the probability $Pr(\{\mathbf{X},\mathbf{Y}\}\in\Psi_\epsilon^c) = O(k/M)$.

**For $\{X,Y\}\in\Psi_\epsilon$** Now note that $(\mathbf{e}_c(X)+\mathbf{e}_t(X))^q$ will contain terms of the form $(\mathbf{e}_c(X))^m(\mathbf{e}_t(X))^{q-m}$. For $m < q$, the term $(\mathbf{e}_c(X))^m(\mathbf{e}_t(X))^{q-m}$ will be a sum of terms of the form $(k/M)^{(m+u)}\mathbf{P}^{-(m+v)}(X)$ for arbitrary $v < q-m$ with $u-v >= 2/d$. By (105), the covariance term $Cov[(\mathbf{e}_c(X))^m(\mathbf{e}_t(X))^{q-m}, (\mathbf{e}_c(Y))^n]$ will be therefore be $O(k_M^{2/d}/M)$ if either $m < q$ or $n < r$.

On the other hand, if $m = q$ and $n = r$, $Cov[(\mathbf{e}_c(X))^q, (\mathbf{e}_c(Y))^r] = 1_{\{q=1,r=1\}}O(1/M) + O(1/M^2)$ by noting that the error $\mathbf{e}_c(X) = f(X)\mathbf{E}_1(X)$ and subsequently invoking (106). Therefore

$$I = \mathbb{E}[1_{\mathbf{\Delta}_\epsilon^c(\mathbf{X},\mathbf{Y})}\sharp(\mathbf{X},\mathbf{Y})]$$
$$= \mathbb{E}\left[1_{\mathbf{\Delta}_\epsilon^c(\mathbf{X},\mathbf{Y})}\left(1_{\{q=1,r=1\}}O(1/M) + O(k_M^{2/d}/M) + O(1/M^2)\right)\right]$$
$$= 1_{\{q=1,r=1\}}O(1/M) + O(k_M^{2/d}/M) + O(1/M^2),$$

where the last step follows from the fact that probability $Pr(\{\mathbf{X},\mathbf{Y}\}\in\Psi_\epsilon) = 1 - O(k/M) = O(1)$. $\qquad\square$

## B.5 Specific cases

We now focus on evaluating the specific cases

$$\mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^2(\mathbf{X})\big]$$

and

$$Cov\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k(\mathbf{Y})\big],$$

for $k > 2$.

### B.5.1 Evaluation of $\mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^2(\mathbf{X})\big]$

$\mathbf{P}(X)$ has a beta distribution with parameters $k, M - k + 1$. Therefore for $k > 2$

$$
\begin{aligned}
\mathbb{E}[\mathbf{E}_\beta^2(X)] &= \mathbb{E}\left[k_M^{2\beta}(\mathbf{P}^{-\beta}(X) - \mathbb{E}[\mathbf{P}^{-\beta}(X)])^2\right] \\
&= k_M^{2\beta}\mathbb{E}[\mathbf{P}^{-2\beta}(X)] - \big(\mathbb{E}[\mathbf{P}^{-\beta}(X)]\big)^2 \\
&= k_M^{2\beta}\left(\frac{\Gamma(k-2\beta)\Gamma(M+1)}{\Gamma(k)\Gamma(M+1-2\beta)} - \left(\frac{\Gamma(k-\beta)\Gamma(M+1)}{\Gamma(k)\Gamma(M+1-\beta)}\right)^2\right) \\
&= O(1/k) \quad (108)
\end{aligned}
$$

where the last step follows by noting that for any $a > 0$,

$$\frac{\Gamma(x)}{\Gamma(x+a)} = x^{-a}(1 + o(1/x)).$$

From ( 103),

$$\mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^2(\mathbf{X})\big] = \mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_o^2(\mathbf{X})\big] + O(\mathcal{C}(k)). \quad (109)$$

Note that $\mathbf{e}_o^2(X) = (\mathbf{e}_c(X)+\mathbf{e}_t(X)+\mathbf{e}_r(X))^2$ is a sum of terms of the form $(\mathbf{e}_c(X))^{2-l-m}(\mathbf{e}_t(X))^l(\mathbf{e}_r(X))^m$. Also,

$$
\begin{aligned}
\mathbb{E}[\mathbf{e}_c^2(X)] &= f^2(X)\mathbb{E}\left[k_M^2(\mathbf{P}^{-1}(X) - \mathbb{E}[\mathbf{P}^{-1}(X)])^2\right] \\
&= f^2(X)k_M^2\mathbb{E}[\mathbf{P}^{-2}(X)] - \big(\mathbb{E}[\mathbf{P}^{-1}(X)]\big)^2 \\
&= f^2(X)k_M^{2\beta}\left(\frac{\Gamma(k-2)\Gamma(M+1)}{\Gamma(k)\Gamma(M+1-2)} - \left(\frac{\Gamma(k-1)\Gamma(M+1)}{\Gamma(k)\Gamma(M)}\right)^2\right) \\
&= \frac{1}{k} + o\left(\frac{1}{k}\right). \quad (110)
\end{aligned}
$$

Using (108), identical to the derivation of (100) and (101), it is clear that if $l + m > 0$, $\mathbb{E}[(\mathbf{e}_c(X))^{2-l-m}(\mathbf{e}_t(X))^l(\mathbf{e}_r(X))^m] = o(k^{-1}) + o(1/M) + O(\mathcal{C}(k))$. This implies that

$$
\begin{aligned}
\mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^2(\mathbf{X})\big] &= \mathbb{E}\big[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_o^2(\mathbf{X})\big] + O(\mathcal{C}(k)) \\
&= f^2(X)\left(\frac{1}{k}\right) + o\left(\frac{1}{k}\right). \quad (111)
\end{aligned}
$$

## B.5.2 Evaluation of $Cov\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k(\mathbf{X}), 1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k(\mathbf{Y})\right]$

We separately analyze disjoint balls and intersecting balls as follows:

$$
\begin{aligned}
&Cov\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k(\mathbf{X}), 1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k(\mathbf{Y})\right] \\
&= \mathbb{E}\left[\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e}_k(\mathbf{X})\mathbf{e}_k(\mathbf{Y})\right]\right] \\
&= \mathbb{E}\left[\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e}_o(\mathbf{X})\mathbf{e}_o(\mathbf{Y})\right]\right] + O(\mathcal{C}(k)) \\
&= \mathbb{E}\left[\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})(\mathbf{e}_c(\mathbf{X}) + \mathbf{e}_t(\mathbf{X}) + \mathbf{e}_r(\mathbf{X}))(\mathbf{e}_c(\mathbf{Y}) + \mathbf{e}_t(\mathbf{Y}) + \mathbf{e}_r(\mathbf{Y}))\right]\right] + O(\mathcal{C}(k)) \\
&= \mathbb{E}\left[\left[1_{\{\mathbf{X} \in \mathcal{S}'\}}1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})(\mathbf{e}_c(\mathbf{X}) + \mathbf{e}_t(\mathbf{X}))(\mathbf{e}_c(\mathbf{Y}) + \mathbf{e}_t(\mathbf{Y}))\right]\right] + O(\mathcal{C}(k)) + o(1/M) \\
&= \mathbb{E}\left[\mathbf{1}_{\boldsymbol{\Delta}_\epsilon{}^c(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbb{E}_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))(\mathbf{e}_c(Y) + \mathbf{e}_t(Y))]\right] \\
&\quad + \mathbb{E}\left[\mathbf{1}_{\boldsymbol{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbb{E}_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))(\mathbf{e}_c(Y) + \mathbf{e}_t(Y))]\right] \\
&\quad + O(\mathcal{C}(k)) + o(1/M) \\
&= I + II + O(\mathcal{C}(k)) + o(1/M).
\end{aligned}
$$

**For $\{X, Y\} \in \Psi_\epsilon$**

$$
\mathbb{E}[(\mathbf{e}_c(X))(\mathbf{e}_c(Y))] = Cov[(\mathbf{e}_c(X)), (\mathbf{e}_c(Y))] = \frac{-f(X)f(Y)}{M} + O(1/M^2)
$$

by noting that the error $\mathbf{e}_c(X) = \mathbf{E}_1(X)/f(X)$ and subsequently invoking (106) in conjunction with the condition $k > 2$. Similarly, using (89), (90) and (106),

$$
\mathbb{E}[(\mathbf{e}_c(X))(\mathbf{e}_t(Y))] = O(k_M^{2/d}/M) + O(1/M^2),
$$

$$
\mathbb{E}[(\mathbf{e}_t(X))(\mathbf{e}_c(Y))] = O(k_M^{2/d}/M) + O(1/M^2),
$$

$$
\mathbb{E}[(\mathbf{e}_t(X))(\mathbf{e}_t(Y))] = O(k_M^{4/d}/M) + O(1/M^2).
$$

This implies that

$$
\begin{aligned}
I &= \mathbb{E}[\mathbf{1}_{\boldsymbol{\Delta}_\epsilon^c(\mathbf{X},\mathbf{Y})}\mathbb{E}_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))(\mathbf{e}_c(Y) + \mathbf{e}_t(Y))]] \\
&= \mathbb{E}\left[\mathbf{1}_{\boldsymbol{\Delta}_\epsilon^c(\mathbf{X},\mathbf{Y})}\left(-f(X)f(Y)(1/M) + O(k_M^{2/d}/M) + O(1/M^2)\right)\right] \\
&= \mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S}'\}}1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})(f(\mathbf{X})f(\mathbf{Y}))]\left(-1/M + O(k_M^{2/d}/M) + O(1/M^2)\right) \\
&= -\mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})f(\mathbf{X})]\mathbb{E}[1_{\{\mathbf{Y} \in \mathcal{S}'\}}\gamma_2(\mathbf{Y})f(\mathbf{Y})]\frac{1}{M} + O(k_M^{2/d}/M) + O(1/M^2). \quad (112)
\end{aligned}
$$

where the last but one step follows from the fact that probability $Pr(\{\mathbf{X}, \mathbf{Y}\} \in \Psi_\epsilon) = 1 - O(k/M) = O(1)$.

**For** $\{X, Y\} \in \Psi_\epsilon^c$ First observe that by Cauchy Schwarz, and by (108) $|\mathbb{E}[\mathbf{E}_t(X)\mathbf{E}_u(X)]| \leq \sqrt{\mathbb{E}[\mathbf{E}_t^2(X)]\mathbb{E}[\mathbf{E}_u^2(X)]} = O(1/k)$. This implies that

$$\mathbb{E}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))(\mathbf{e}_c(Y) + \mathbf{e}_t(Y))] = \mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_c(Y)] + O(k_M^{2/d}/k). \tag{113}$$

In subsection B.7, we will show Lemma B.5, which states that

$$\mathbb{E}[\mathbf{1}_{\Delta_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e_c}(\mathbf{X})\mathbf{e_c}(\mathbf{Y})]$$
$$= \mathbb{E}[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})]\left(\frac{1}{M} + o\left(\frac{1}{M}\right)\right)$$

This implies that

$$II = \mathbb{E}[\mathbf{1}_{\Delta_\epsilon(\mathbf{X},\mathbf{Y})}\mathbb{E}_{\{\mathbf{X},\mathbf{Y}\}}[(\mathbf{e}_c(X) + \mathbf{e}_t(X))(\mathbf{e}_c(Y) + \mathbf{e}_t(Y))]]$$
$$= \mathbb{E}[\mathbf{1}_{\Delta_\epsilon(\mathbf{X},\mathbf{Y})}\mathbb{E}_{\{\mathbf{X},\mathbf{Y}\}}[\mathbf{e}_c(X)\mathbf{e}_c(Y)] + O(k_M^{2/d}/k)]$$
$$= \mathbb{E}[\mathbf{1}_{\Delta_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e_c}(\mathbf{X})\mathbf{e_c}(\mathbf{Y})] + \mathbb{E}\left[\mathbf{1}_{\Delta_\epsilon(\mathbf{X},\mathbf{Y})}\left(O(k_M^{2/d}/k)\right)\right]$$
$$= \mathbb{E}[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})/f^2(\mathbf{X})]\left(\frac{1}{M} + O(k_M^{2/d}/M) + o\left(\frac{1}{M}\right)\right) \tag{114}$$

where the last step follows from recognizing that $Pr(\{\mathbf{X}, \mathbf{Y}\} \in \Psi_\epsilon^c) = O(k/M)$ and $O(k/M) \times 1/k = O(1/M)$. This implies that

$$Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k(\mathbf{Y})\right]$$
$$= I + II + O(\mathcal{C}(k)) + o(1/M)$$
$$= Cov[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})/f(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})/f(\mathbf{Y})]\left(\frac{1}{M}\right) + o(1/M) + O(\mathcal{C}(k)). \tag{115}$$

## B.6   Summary

Noting that $\delta > 2/3$, the equations (83), (B.2), (111), (115) imply that for positive integers $q, r < k$,

$$\mathbb{E}\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})\mathbf{e}_k^q(\mathbf{X})\right] = 1_{\{q=2\}}\mathbb{E}\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma(\mathbf{X})f^2(\mathbf{X})\right]\left(\frac{1}{k}\right) + o\left(\frac{1}{k}\right) + O(\mathcal{C}(k)), \tag{116}$$

$$Cov\left[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\mathbf{e}_k^q(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})\mathbf{e}_k^r(\mathbf{Y})\right]$$
$$= 1_{\{q,r=1\}}Cov[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})f(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})f(\mathbf{Y})]\left(\frac{1}{M} + o(1/M)\right)$$
$$+ 1_{\{q+r>2\}}\left(O\left(\frac{1}{k^{((q+r)\delta/2-1)}M}\right) + O(k_M^{2/d}/M) + O(1/M^2)\right) + O(\mathcal{C}(k)). \tag{117}$$

67

## B.7 Evaluation of $\mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_c(Y)]$ for $\{X, Y\} \in \Psi_\epsilon^c$

For $\{X, Y\} \in \Psi_\epsilon^c$, it will be shown that the cross-correlations $\mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_c(Y)]$ of the coverage density estimator and an oracle uniform kernel density estimator (defined below) are identical up to leading terms (without explicitly evaluating the cross-correlation between the coverage density estimates) and then derive the correlation of the oracle density estimator to obtain corresponding results for the coverage estimate.

**Oracle $\epsilon$ ball density estimate** In order to estimate cross moments for the $k$-NN density estimator, the $\epsilon$ *ball* density estimator is introduced. The $\epsilon$-ball density estimator is a kernel density estimator that uses a uniform kernel with bandwidth which depends on the unknown density $f$. Let the volume of the kernel be $V_\epsilon(X)$ and the corresponding kernel region be $S_\epsilon(X) = \{Y \in \mathcal{S} : c_d \|X - Y\|^d \leq V_\epsilon(X)\}$. The volume is chosen such that the coverage $Q_\epsilon(X) = \int_{S_\epsilon(X)} f(z)dz$ is set to $(1 + p_k)k/M$. Let $\mathbf{l}_\epsilon(X)$ denote the number of points among $\{\mathbf{X}_1, .., \mathbf{X}_M\}$ falling in $S_\epsilon(X)$: $\mathbf{l}_\epsilon(\mathbf{X}) = \Sigma_{i=1}^M 1_{\mathbf{X}_i \in S_\epsilon(X)}$. The $\epsilon$ *ball* density estimator is defined as

$$\hat{\mathbf{f}}_\epsilon(X) = \frac{\mathbf{l}_\epsilon(\mathbf{X})}{MV_\epsilon(X)}. \tag{118}$$

Also define the error $\mathbf{e}_\epsilon(X)$ as $\mathbf{e}_\epsilon(X) = \hat{\mathbf{f}}_\epsilon(X) - \mathbb{E}[\hat{\mathbf{f}}_\epsilon(X)]$. It is then possible to prove the following lemma using results on the volumes of intersections of hyper spheres (refer Appendix A for details).

**Lemma B.3.** *Let $\gamma_1(X)$, $\gamma_2(X)$ be arbitrary continuous functions. Let $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ denote $M + 2$ i.i.d realizations of the density $f$. Then,*

$$\mathbb{E}\left[\mathbf{1}_{\boldsymbol{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\mathbf{e}_\epsilon(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e}_\epsilon(\mathbf{Y})\right]$$
$$= \mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})]\left(\frac{1}{M} + o\left(\frac{1}{M}\right)\right).$$

Next, the cross-correlations of the coverage density estimator and the $\epsilon$ ball density estimator are shown to be asymptotically equal. In particular,

**Lemma B.4.**
$$\mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_c(Y)] = \mathbb{E}[\mathbf{e}_\epsilon(X)\mathbf{e}_\epsilon(Y)] + o(1/k).$$

*Proof.* We begin by establishing the conditional density and expectation of $\hat{\mathbf{f}}_\epsilon(X)$ given $\hat{\mathbf{f}}_\mathbf{c}(X)$. We drop the dependence on $X$ and denote $\mathbf{l}_\epsilon = \Sigma_{i=1}^M 1_{\{X_i \in S_\epsilon(X)\}}$, the $k$-NN coverage by $\mathbf{P}$ and the $\epsilon$ ball coverage $Q_\epsilon(X)$ by $Q$. Let $\mathbf{q} = Q/\mathbf{P}$ and $\mathbf{r} = (Q - \mathbf{P})/(1 - \mathbf{P})$. The following expressions for conditional densities and expectations are derived in [33]

$$\mathbf{Pr}\{\mathbf{l}_\epsilon = l|\mathbf{P}; \mathbf{P} > Q\}$$
$$= \begin{cases} \binom{k-1}{l}\mathbf{q}^l(1 - \mathbf{q})^{k-1-l} & l = 0, 1, \ldots, k - 1 \\ 0 & l = k, k + 1, \ldots, M \end{cases}$$

68

$$\mathbf{Pr}\{\mathbf{l}_\epsilon = l | \mathbf{P}; \mathbf{P} \leq Q\}$$
$$= \begin{cases} 0 & l = 0, 1, \ldots, k-1 \\ \binom{M-k}{l-k}\mathbf{r}^{l-k}(1-\mathbf{r})^{M-l} & l = k, k+1, \ldots, M \end{cases}$$

which implies

$$\mathbb{E}[\mathbf{l}_\epsilon = l | \mathbf{P}; \mathbf{P} > Q] = (k-1)Q/\mathbf{P}$$
$$\mathbb{E}[\mathbf{l}_\epsilon = l | \mathbf{P}; \mathbf{P} \leq Q] = \left(\frac{1-Q}{1-\mathbf{P}}\right)(k-M) + M$$

Using the above expressions for conditional expectations, the following marginal expectation are obtained. Denote the density of the coverage $\mathbf{P}$ by $f_{k,M}(p)$. Also let $\hat{\mathbf{P}}$ be the coverage corresponding to the $k-2$ nearest neighbor in a total field of $M-3$ points. Then

$$
\begin{aligned}
\mathbb{E}[\tilde{\mathbf{e}}_c(X)\tilde{\mathbf{e}}_\epsilon(X)] &= \mathbb{E}[\hat{\mathbf{f}}_\epsilon(X)\hat{\mathbf{f}}_\mathbf{c}(X)] - \mathbb{E}[\hat{\mathbf{f}}_\mathbf{c}(X)]\mathbb{E}[\hat{\mathbf{f}}_\epsilon(X)] \\
&= \mathbb{E}\left[\left(\left(\frac{1-Q}{\mathbf{P}(1-\mathbf{P})}\right)(k-M) + M/\mathbf{P}\right)1_{\mathbf{P}\leq Q}\right] \\
&\quad + \frac{f^2(X)(k-1)}{kM}\mathbb{E}\left[((k-1)Q/\mathbf{P}^2)\,1_{\mathbf{P}>Q}\right] - \frac{f^2(X)}{k}MQ. \\
&= \frac{f^2(X)}{k}\frac{(M-1)(M-2)}{(k-2)(M-k)} \times \\
&\quad \mathbb{E}[(1-Q\hat{\mathbf{P}})(k-M) + M\hat{\mathbf{P}}(1-\hat{\mathbf{P}})] - \frac{f^2(X)}{k}MQ \\
&\quad + \mathbb{E}[((k-1)Q(1-\hat{\mathbf{P}}) - (1-Q\hat{\mathbf{P}})(k-M) + M\hat{\mathbf{P}}(1-\hat{\mathbf{P}}))(1_{\hat{\mathbf{P}}>Q})] \\
&= C \times (I - II + III).
\end{aligned}
$$

It can be shown that $C \times (I-II) = \frac{f^2(X)}{k}(1-Q)$ using the fact that $\hat{\mathbf{P}}$ has a beta distribution. Note that from the definition of $Q = ((1+p_k)(k-1)/M)$, from the concentration inequality we have that $\mathbb{E}[1_{\hat{\mathbf{P}}>Q}] = \mathcal{C}(M)$. The remainder $(C \times III)$ can be simplified and bounded using the Cauchy-Schwarz inequality and the concentration inequality to show $C \times III = o(1/M)$.

Therefore,

$$
\begin{aligned}
\mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_\epsilon(X)] &= \frac{f^2(X)}{k}(1-Q) + \mathcal{C}(M). \\
&= \frac{f^2(X)}{k} - \frac{f^2(X)}{M} + o\left(\frac{1}{M}\right) \\
&= f^2(X)\left(\frac{1}{k} + o\left(\frac{1}{k}\right)\right). \quad (119)
\end{aligned}
$$

Now denote $\mathbf{E}(X) = (\mathbf{e}_c(X) - \mathbf{e}_\epsilon(X))$. Note that $\mathbb{E}[\mathbf{E}^2(X)] = \mathbb{E}[\mathbf{e}_c(X)^2] - 2E[\mathbf{e}_c(X)\mathbf{e}_\epsilon(X)] + \mathbb{E}[\mathbf{e}_\epsilon(X)^2]$. Since $E[\mathbf{e}_c(X)^2] = f^2(X)\frac{1}{k} + o(1/k)$ and $E[\mathbf{e}_\epsilon(X)^2] = f^2(X)(1/k + o(1/k))$ it follows from (119) that $\mathbb{E}[E(X)] = o(1/k)$. This result means $\mathbf{e}_c(X)$ and $\mathbf{e}_\epsilon(X)$ are almost

perfectly correlated. Next express the covariance between the coverage density estimates in terms of the covariance between the $\epsilon$ ball estimates as follows:

$$
\begin{aligned}
&\mathbb{E}[\mathbf{e}_c(X)\mathbf{e}_c(Y)] \\
&= \mathbb{E}[(\mathbf{e}_\epsilon(X) + \mathbf{E}(X))(\mathbf{e}_\epsilon(Y) + \mathbf{E}(Y))] \\
&= \mathbb{E}[\mathbf{e}_\epsilon(X)\mathbf{e}_\epsilon(Y)] + \mathbb{E}[\mathbf{e}_\epsilon(X)(\mathbf{E}(Y))] \\
&\quad + \mathbb{E}[\mathbf{e}_\epsilon(Y)(\mathbf{E}(X))] + \mathbb{E}[(\mathbf{E}(X))(\mathbf{E}(Y))] \\
&= I + II + III + IV.
\end{aligned}
$$

Using Cauchy-Schwarz, a bound on each of the terms $II$, $III$ and $IV$ is obtained in terms of $\mathbb{E}[\mathbf{E}(X)]$: $|II| \leq \sqrt{\mathbb{E}[\mathbf{E}(Y)]\mathbb{E}[\mathbf{e}_\epsilon{}^2(X)]}$, $|III| \leq \sqrt{\mathbb{E}[\mathbf{E}(X)]\mathbb{E}[\mathbf{e}_\epsilon{}^2(Y)]}$ and $|IV| \leq \sqrt{\mathbb{E}[\mathbf{E}(X)]\mathbb{E}[\mathbf{E}(Y)]}$. Note that the above application of Cauchy-Schwarz *decouples* the problem of joint expectation of density estimates located at two *different* points $X$ and $Y$ to a problem of estimating the error $\mathbf{E}$ between two different density estimates at the *same* point(s). Therefore all the three terms $II$, $III$ and $IV$ are $o(1/k)$. This concludes the proof of Lemma B.4. □

For Lemma B.4 to be useful, $\mathbb{E}[\mathbf{e}_\epsilon(X)\mathbf{e}_\epsilon(Y)]$ must be orders of magnitude larger than the error $o(1/k)$, which is indeed the case for $\{X, Y\} \in \Psi_\epsilon{}^c$ since $\mathbb{E}[\mathbf{e}_\epsilon(X)\mathbf{e}_\epsilon(Y)] = O(1/k)$ (Lemma A.2, Appendix .1) for such $X$ and $Y$. This lemma can be used along with previously established results on co-variance of $\epsilon$-ball density estimates (Lemma B.3) to obtain the following result:

**Lemma B.5.** *Let* $\gamma_1(X)$, $\gamma_2(X)$ *be arbitrary continuous functions. Let* $\mathbf{X}_1, .., \mathbf{X}_M, \mathbf{X}, \mathbf{Y}$ *denote* $M + 2$ *i.i.d realizations of the density* $f$. *Then,*

$$
\begin{aligned}
&\mathbb{E}[\mathbf{1}_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e_c}(\mathbf{X})\mathbf{e_c}(\mathbf{Y})] \\
&= \mathbb{E}[\mathbf{1}_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})]\left(\frac{1}{M} + o\left(\frac{1}{M}\right)\right)
\end{aligned}
$$

*Proof.*

$$
\begin{aligned}
&\mathbb{E}[\mathbf{1}_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbb{E}_{\mathbf{X},\mathbf{Y}}[\mathbf{e_c}(X)\mathbf{e_c}(Y)]] \\
&= \mathbb{E}[\mathbf{1}_{\mathbf{\Delta}_\epsilon(\mathbf{X},\mathbf{Y})}\gamma_1(\mathbf{X})\gamma_2(\mathbf{Y})\mathbf{e}_\epsilon(\mathbf{X})\mathbf{e}_\epsilon(\mathbf{Y})] + o(1/k) \\
&= \mathbb{E}[\mathbf{1}_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})\gamma_2(\mathbf{X})f^2(\mathbf{X})]\left(\frac{1}{M} + o\left(\frac{1}{M}\right)\right).
\end{aligned}
$$

In the second to last step, $o(1/M)$ is obtained for the second term by recognizing that $Pr(\{\mathbf{X}, \mathbf{Y}\} \in \Psi_\epsilon^c) = O(k/M)$ and $O(k/M) \times o(1/k) = o(1/M)$. □

# C  Boundary correction for density estimates

In the previous section, moment results were established for the standard $k$-NN density estimate $\hat{\mathbf{f}}_k(X)$ for points $X$ in any deterministic set $\mathcal{S}'$ with respect to the samples $\mathcal{X}_M =$

$\{\mathbf{X}_{N+1}, .., \mathbf{X}_{N+M}\}$ satisfying the condition $Pr(\mathbf{X} \notin \mathcal{S}') = o(1)$ and $\mathcal{S}' \subset \mathcal{S}_I$, where $\mathbf{X}$ is an realization from density $f$. In this section, these moment results are extended to boundary corrected $k$-NN density estimate $\tilde{\mathbf{f}}_k(X)$ for all $X \in \mathcal{S}$ as follows.

Specify the set $\mathcal{S}'$ to be $\mathcal{S}' = \mathcal{S}_I$ as defined in (2). Exclusively using the set $\mathcal{X}_N = \{\mathbf{X}_1, .., \mathbf{X}_N\}$, a set of interior points $\mathcal{I}_N \subset \mathcal{X}_N$ are determined such that $\mathcal{I}_N \subset \mathcal{S}'$ with high probability $1 - O(N\mathcal{C}(k))$. Define the set of boundary points $\mathcal{B}_N = \mathcal{X}_N - \mathcal{I}_N$. For points $X \in \mathcal{I}_N$, the boundary corrected $k$-NN density estimate $\tilde{\mathbf{f}}_k(X)$ is defined to be the standard $k$-NN estimate $\hat{\mathbf{f}}_k(X)$, and we invoke the moment properties of the standard $k$-NN density estimate $\hat{\mathbf{f}}_k(X)$ derived in the previous section. For points $X \in \mathcal{B}_N$, the density estimate $\tilde{\mathbf{f}}_k(X)$ is defined as $\hat{\mathbf{f}}_k(Y_n)$ for points $Y_n \in \mathcal{I}_N$, and we invoke the moment properties of the standard $k$-NN density estimate $\hat{\mathbf{f}}_k(X)$ derived in the previous section.

## C.1   Bias in the $k$-NN density estimator near boundary

If a probability density function has bounded support, the $k$-NN balls centered at points close to the boundary are often truncated at the boundary. Let

$$\alpha_k(X) = \frac{\int_{\mathbf{S}_k(X) \cap \mathcal{S}} dZ}{\int_{\mathbf{S}_k(X)} dZ}$$

be the fraction of the volume of the $k$-NN ball inside the boundary of the support. Also define $\mathbf{V}_{k,M}(X)$ to be the $k$-NN ball volume in a sample of size $M$. For interior points $X \in \mathcal{S}'$, $\alpha_k(X) = 1$, while for boundary points $X \in \mathcal{S} - \mathcal{S}'$, $\alpha_k(X)$ is closer to 0 when the points are closer to the boundary. For boundary points we then have

$$\mathbb{E}[\hat{\mathbf{f}}_k(X)] - f(X) = (1 - \alpha_k(X))f(X) + o(1). \tag{120}$$

Therefore the bias is much higher at the boundary of the support ($O(1)$) as compared to its interior ($O((k/M)^{2/d})$) (81). Furthermore, the bias at the support boundary does not decay to 0 as $k/M \to 0$.

In the next section, we detect interior points $\mathcal{I}_N$ which lie in $\mathcal{S}'$ with high probability $O(N\mathcal{C}(k))$. The results on bias, variance and cross-moments derived in the previous Appendix for points $X \in \mathcal{S}'$ therefore carry over to the points $\mathcal{I}_N$. A density estimate at points $\mathcal{B}_N$ is then proposed that will reduce the bias of density estimates close to the boundary.

## C.2   Boundary point detection

Define $V_{k,M}(X) := \frac{k}{M\alpha_k(X)f(X)}$. Let $p(k, M)$ be any positive function satisfying $p(k, M) = \Theta((k/M)^{2/d}) + (\sqrt{6}/k^{\delta/2})$. From the concentration inequality (58) and Taylor series expansion of the coverage function (70), for small values of $k/M$, we have

$$1 - Pr\left(\left|\frac{\mathbf{V}_{k,M}(X)}{V_{k,M}(X)} - 1\right| \le p(k, M)\right) = O(\mathcal{C}(k)).$$

To determine $\mathcal{I}_\mathbb{N}$ and $\mathcal{B}_\mathbb{N}$, we first construct a $K$-NN graph on the samples $\mathcal{X}_N$ where $K = \lfloor k \times (N/M) \rfloor$. For any $X \in \mathcal{X}_\mathbb{N}$, from the concentration inequality (58)

$$1 - Pr\left(\left|\frac{\mathbf{V}_{K,N}(X)}{V_{K,N}(X)} - 1\right| \leq p(K, N)\right) = O(\mathcal{C}(K)) = O(\mathcal{C}(k)), \tag{121}$$

where $\mathcal{C}(K) = O(\mathcal{C}(k))$ because by $(\mathcal{A}.0)$, $K = \theta(k)$. This implies that, with high probability, the radius of the $K$-NN ball at $X$ concentrates around $(V_{K,N}(X)/c_d)^{1/d}$. By this concentration inequality (121), this choice of $K$ guarantees that the size of the $k$-NN ball in the partitioned sample is the same as the the size of the $K$-NN ball in the pooled sample with high probability $1 - \mathcal{C}(k)$. By the union bound and (121), the probability that

$$\left|\frac{\mathbf{V}_{K,N}(X)}{V_{K,N}(X)} - 1\right| \leq p(K, N)$$

is satisfied by every $X_i \in \mathcal{X}_N$ is lower bounded by $1 - O(N\mathcal{C}(k))$.

Using the $K$-NN graph, for each sample $\mathbf{X} \in \mathcal{X}_\mathbb{N}$, we compute the number of points in $\mathcal{X}_\mathbb{N}$ that have $\mathbf{X}$ as a $l$-th nearest neighbor ($l$-NN), $l = \{1, \ldots, K\}$. Denote this count as $count(\mathbf{X})$. Let $Y$ be the $l$-nearest neighbor of $X$, $l = \{1, \ldots, K\}$. Then $Y$ can be represented as $Y = X + R_K(X)u$ where $u$ is an arbitrary vector with $||u|| \leq 1$.

For $X$ to be one of the $K$-NN of $Y$ it is necessary that $R_K(Y) \geq ||Y - X||$ or equivalently, $R_K(Y)/R_K(X) \geq ||u||$. Using the concentration inequality (121) for $R_K(X)$ and $R_K(Y)$, a sufficient condition for this is

$$\frac{\alpha_K(X)f(X)}{\alpha_K(Y)f(Y)}(1 - 2p(K, N)) \geq ||u||. \tag{122}$$

Because $f$ is differentiable and has a finite support, $f$ is Lipschitz continuous. Denote the Lipschitz constant by $\mathbb{L}$. Then, we have $|f(Y) - f(X)| \leq \mathbb{L}(K/c_d N \epsilon_0)^{1/d}$. Define $q(K, N) = (\mathbb{L}/\epsilon_0)(K/c_d N \epsilon_0)^{1/d} + 2\sqrt{6}/k^{\delta/2}$. Then (122) is satisfied if

$$\frac{\alpha_K(X)}{\alpha_K(Y)}(1 - q(K, N)) \geq ||u||.$$

For points $X \in \mathcal{S}'$, $\alpha_K(X) = 1$ with probability $1 - \mathcal{C}(k)$. This implies that $X$ will be one of the $K$-NN of $Y$ if $||u|| \leq 1 - q(K, N)$. This implies that, with probability $1 - O(N\mathcal{C}(k))$, $count(\mathbf{X}) \geq K(1-q(K, N))$ whenever $X \in \mathcal{S}'$. On the other hand, for $X \in \mathcal{S} - \mathcal{S}'$, $\alpha_K(X) < 1$ with probability $1 - \mathcal{C}(k)$. It is also clear that for small values of $K/N$, $\alpha_K(X) < \alpha_K(Y)$ for at least $K/2$ $l$-NN $Y$ of $X$. This then implies that $count(\mathbf{X}) < K(1-q(K, N))$ for $X \in \mathcal{S} - \mathcal{S}'$ with probability $1 - O(N\mathcal{C}(k))$. We therefore can apply the threshold $K(1 - q(K, N))$ to detect interior points $\mathcal{I}_\mathbb{N} = \mathcal{X}_\mathbb{N} \cap \mathcal{S}'$ and boundary points $\mathcal{B}_N = \mathcal{X}_\mathbb{N} - \mathcal{I}_\mathbb{N} = \mathcal{X}_\mathbb{N} \cap (\mathcal{S} - \mathcal{S}')$ with high probability $1 - O(N\mathcal{C}(k))$. Algorithm 1, shown below, codifies this into a precise procedure.

**Algorithm 1** Detect boundary points $\mathcal{B}_N$

---

   1. Construct $K$-NN tree on $\mathfrak{X}_N$

   2. Compute $count(\mathbf{X})$ for each $\mathbf{X} \in \mathfrak{X}_N$

   3. Detect boundary points $\mathcal{B}_N$:

  **for** each $\mathbf{X} \in \mathfrak{X}_N$ **do**

    **if** $count(\mathbf{X}) < (1 - q(K, N))K$ **then**

      $\mathcal{B}_N \leftarrow \mathbf{X}$

    **else**

      $\mathcal{I}_N \leftarrow \mathbf{X}$

    **end if**

  **end for**

---

## C.3   Boundary corrected density estimator

Here the boundary corrected $k$-NN density estimator is defined and its asymptotic rates are computed. The proposed density estimator corrects the $k$-NN ball volumes for points that are close to the boundary. To estimate the density at a boundary point $\mathbf{X} \in \mathcal{B}_N$, we find a point $\mathbf{Y} \in \mathcal{I}_N$ that is close to $\mathbf{X}$. Because of the proximity of $\mathbf{X}$ and $\mathbf{Y}$, $f(\mathbf{X}) \approx f(\mathbf{Y})$. We can then estimate the density at $\mathbf{Y}$ instead and use this as an estimate of $f(\mathbf{Y})$. This informal argument is made more precise in what follows.

Consider the corrected density estimator $\tilde{\mathbf{f}}_k$ defined in (3). This estimator has bias of order $O((k/M)^{1/d})$, which can be shown as follows. Let $\mathbf{X}$ denote $\mathbf{X}_i$ for some fixed $i \in \{1, .., N\}$. Also, let $\mathbf{X}_{-1} = \arg\min_{x \in \mathcal{S}'} d(x, \mathbf{X})$.

Given $\mathfrak{X}_N$, if $X \in \mathcal{I}_N$, then by (81),

$$\mathbb{E}[\tilde{\mathbf{f}}_k(X)] \;\; = \;\; \mathbb{E}[\hat{\mathbf{f}}_k(X)] = f(X) + O((k/M)^{2/d}) + O(\mathcal{C}(k)).$$

Next consider the alternative case $X \in \mathcal{B}_N$. Let $X_n \in \mathcal{I}_N$ be the closest interior point to $X$. Define $h = X - X_n$. $h$ can be rewritten as $h = h_1 + h_2$, where $h_1 = X - X_{-1}$ and $h_2 = X_{-1} - X_n$. Since $X \in \mathcal{B}_N$ implies that $X \in \mathcal{S} - \mathcal{S}'$ with probability $1 - O(N\mathcal{C}(k))$, consequently $||h_1|| = ||X - X_{-1}|| = O((k/M)^{1/d})$ with probability $1 - O(N\mathcal{C}(k))$.

Again with probability $1 - O(N\mathcal{C}(k))$, $X_n \in \mathcal{S}'$. Let $\mathcal{C}_N = \cup_{Y \in \mathcal{S}'} \arg\min_{x \in \mathcal{I}_N} d(x, Y)$. By construction of $\mathcal{C}_N$, $X_n \in \mathcal{C}_N$. Consequently, by (121), $||h_2|| = ||X_{-1} - X_n|| = O((1/N)^{1/d}) = o((k/M)^{1/d})$.

Because $||h_1|| = ||X - X_{-1}|| = O((k/M)^{1/d})$ and $||h_2|| = ||X_{-1} - X_n|| = o((k/M)^{1/d})$ with probability $1 - O(N\mathcal{C}(k))$, consequently with probability $1 - O(N\mathcal{C}(k))$, $||h|| = O((k/M)^{1/d})$. Now,

$$f(X) = f(X_n) + O(||h||).$$

If $X_n$ is located in the interior $\mathcal{S}'$, by (81),

$$\mathbb{E}[\hat{\mathbf{f}}_k(X_n)] \;\; = \;\; f(X_n) + O((k/M)^{2/d}) + O(\mathcal{C}(k)), \tag{123}$$

and therefore

$$
\begin{aligned}
\mathbb{E}[\tilde{\mathbf{f}}_k(X)] &= \mathbb{E}[\hat{f}_k(\mathbf{X}_n)] + O(N\mathcal{C}(k)) \\
&= f(X_n) + O((k/M)^{2/d}) + O(N\mathcal{C}(k)) \\
&= f(X) + O(\|h\|) + O((k/M)^{2/d}) + O(N\mathcal{C}(k)) \\
&= f(X) + O((k/M)^{1/d}) + O(N\mathcal{C}(k)),
\end{aligned}
\tag{124}
$$

where the $O(N\mathcal{C}(k))$ accounts for error in the case of the event that $X_{n(i)} \notin \mathcal{S}'$. This implies that the corrected density estimate has lower bias as compared to the standard $k$-NN density estimate (compare to (81) and (120)). In particular, boundary compensation has reduced the bias of the estimator at points near the boundary from $O(1)$ to $O((k/M)^{1/d}) + O(N\mathcal{C}(k))$.

## C.4 Properties of boundary corrected density estimator

By section C.2, $\mathcal{I}_N \in \mathcal{S}'$ with probability $1 - N\mathcal{C}(k)$. The results on bias, variance and cross-moments of the standard $k$-NN density estimator $\hat{\mathbf{f}}_k$ derived in the previous Appendix for points $X \in \mathcal{S}'$ therefore carry over to the corrected density estimator $\tilde{\mathbf{f}}_k$ for points $\mathcal{I}_N$ with error of order $O(N\mathcal{C}(k))$.

In the definition of the corrected estimator $\tilde{\mathbf{f}}_k$ in (3), $\hat{\mathbf{f}}_k(\mathbf{X}_{n(i)})$ is the standard $k$-NN density estimates and $\mathbf{X}_{n(i)} \in \mathcal{S}'$ . It therefore follows that the variance and other central and cross moments of the corrected density estimator $\tilde{\mathbf{f}}_k$ will continue to decay at the same rate as the standard $k$-NN density estimator in the interior, as given by (116) and (117).

Given these identical rates and that the probability of a point being in the boundary region $\mathcal{S} - \mathcal{S}'$ is $O((k/M)^{1/d}) = o(1)$, the contribution of the boundary region to the overall variance and other cross moments of the boundary corrected density estimator $\tilde{\mathbf{f}}_k$ are asymptotically negligible compared to the contribution from the interior. As a result we can now generalize the results from Appendix A on the central moments and cross moments to include the boundary regions as follows. Denote $\tilde{\mathbf{f}}_k(X) - \mathbb{E}_X[\tilde{\mathbf{f}}_k(X) \mid X]$ by $\mathbf{e}(X)$.

### C.4.1 Central and cross moments

For positive integers $q, r < k$

$$
\mathbb{E}[\gamma(\mathbf{X})\mathbf{e}^q(\mathbf{X})] = 1_{\{q=2\}}\mathbb{E}\big[\gamma(\mathbf{X})f^2(\mathbf{X})\big]\left(\frac{1}{k}\right) + o\left(\frac{1}{k}\right) + O(N\mathcal{C}(k)),
\tag{125}
$$

$$
\begin{aligned}
&Cov[\gamma_1(\mathbf{X})\mathbf{e}^q(\mathbf{X}), \gamma_2(\mathbf{Y})\mathbf{e}^r(\mathbf{Y})] \\
&= 1_{\{q,r=1\}}Cov[1_{\{\mathbf{X}\in\mathcal{S}'\}}\gamma_1(\mathbf{X})f(\mathbf{X}), 1_{\{\mathbf{Y}\in\mathcal{S}'\}}\gamma_2(\mathbf{Y})f(\mathbf{Y})]\left(\frac{1}{M} + o(1/M)\right) \\
&\quad + 1_{\{q+r>2\}}\left(O\left(\frac{1}{k^{((q+r)\delta/2-1)}M}\right) + O(k_M^{2/d}/M) + O(1/M^2)\right) + O(N\mathcal{C}(k)).
\end{aligned}
\tag{126}
$$

Next, we derive the following result on the bias of boundary corrected estimators.

### C.4.2 Bias

For $k > 2$,

$$\mathbb{E}[\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(\mathbf{X}) \mid \mathbf{X}]) - \gamma(f(\mathbf{X})))] = \mathbb{E}\left[\mathbb{E}\left[(\gamma(\tilde{\mathbf{f}}_k(\mathbf{X})) - \gamma(f(\mathbf{X}))) \mid \mathcal{X}_N\right]\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[1_{\{X \in \mathcal{I}_N\}}(\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(X)]) - \gamma(f(X))) \mid \mathcal{X}_N\right]\right] + \mathbb{E}\left[\mathbb{E}\left[1_{\{X \in \mathcal{B}_N\}}(\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(X)]) - \gamma(f(X))) \mid \mathcal{X}_N\right]\right]$$

$$= I + II. \tag{127}$$

From (81), and $Pr(\mathbf{X} \in \mathcal{B}_N) = O((k/M)^{1/d})$, we have

$$I = \mathbb{E}\left[\gamma'(f(\mathbf{X}))h(\mathbf{X})\right]\left(\frac{k}{M}\right)^{2/d} + o\left(\frac{k}{M}\right)^{2/d} + O(N\mathcal{C}(k)). \tag{128}$$

Next, we will now derive $II$.

$$II = \mathbb{E}\left[\mathbb{E}\left[1_{\{X \in \mathcal{B}_N\}}(\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(X)]) - \gamma(f(X))) \mid \mathcal{X}_N\right]\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[1_{\{X \in \mathcal{B}_N\}}(\gamma(f(X_n)) - \gamma(f(X))) + O\left(\frac{k}{M}\right)^{2/d} \mid \mathcal{X}_N\right]\right] + O(N\mathcal{C}(k)) \tag{129}$$

where the last step follows by (123). Let us concentrate on the inner expectation now. By section C.2, we know that with probability $1 - O(N\mathcal{C}(k))$, if $X \in \mathcal{B}_N$, then $X \in \mathcal{S} - \mathcal{S}'$ and if $X_n \in \mathcal{I}_N$, then $X_n \in \mathcal{S}'$. Furthermore, $||X - X_{-1}|| = O(k/M)^{1/d}$ and $||X_{-1} - X_n|| = o(k/M)^{1/d}$ with probability $1 - O(N\mathcal{C}(k))$. This implies that

$$\mathbb{E}\left[1_{\{X \in \mathcal{B}_N\}}(\gamma(f(X_n)) - \gamma(f(X))) + O\left(\frac{k}{M}\right)^{2/d} \mid \mathcal{X}_N\right]$$

$$= \mathbb{E}\left[1_{\{X \in \mathcal{S} - \mathcal{S}'\}}(\gamma(f(X_{-1})) - \gamma(f(X))) \mid \mathcal{X}_N\right] + o\left(\frac{k}{M}\right)^{1/d} + O(N\mathcal{C}(k)).$$

Since $Pr(\mathbf{X} \in \mathcal{S} - \mathcal{S}') = O((k/M)^{1/d})$, this in turn implies that

$$II = \mathbb{E}\left[\mathbb{E}\left[1_{\{X \in \mathcal{B}_N\}}(\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(X)]) - \gamma(f(X))) \mid \mathcal{X}_N\right]\right]$$

$$= \mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S} - \mathcal{S}'\}}(\gamma(f(\mathbf{X}_{-1})) - \gamma(f(\mathbf{X})))] + o\left(\frac{k}{M}\right)^{2/d} + O(N\mathcal{C}(k)). \tag{130}$$

We therefore finally get,

$$\mathbb{E}[\gamma(\mathbb{E}[\tilde{\mathbf{f}}_k(\mathbf{X}) \mid \mathbf{X}]) - \gamma(f(\mathbf{X})))] = I + II$$

$$= \mathbb{E}\left[\gamma'(f(\mathbf{X}))h(\mathbf{X})\right]\left(\frac{k}{M}\right)^{2/d} + \mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S} - \mathcal{S}'\}}(\gamma(f(\mathbf{X}_{-1})) - \gamma(f(\mathbf{X})))] + o\left(\frac{k}{M}\right)^{2/d} + O(N\mathcal{C}(k)). \tag{131}$$

Note that $||\mathbf{X} - \mathbf{X}_{-1}|| = O((k/M)^{1/d})$ with probability $1 - O(N\mathcal{C}(k))$. This therefore implies that

$$c_3 = \mathbb{E}[1_{\{\mathbf{X} \in \mathcal{S} - \mathcal{S}'\}}(\gamma(f(\mathbf{X}_{-1})) - \gamma(f(\mathbf{X})))] = O((k/M)^{1/d}) \times O((k/M)^{1/d}) + O(N\mathcal{C}(k)) = O((k/M)^{2/d}) + O(N\mathcal{C}($$

### C.4.3 Optimality of boundary correction

Comparing (131), (125) and (126) with (81), (116) and (117) respectively, oracle rates of convergence of bias, and central and cross moments for the boundary corrected density estimate are attained. The oracle rates are defined as the rates of MSE convergence attainable by the *oracle* density estimate that knows the boundary of $\mathcal{S}$

$$\tilde{\mathbf{f}}_{k,o} = \frac{k-1}{M\mathbf{V}_{k,o}(X)},$$

where $\mathbf{V}_{k,o}(X)$ is the volume of the region $\mathbf{S}_k(X) \cap \mathcal{S}$. It follows that the boundary compensated BPI estimator is adaptive in the sense that it's asymptotic MSE rate of convergence is identical to that of a $k$-NN plug-in estimator that knows the true boundary. Equivalent corrections exist for the uniform kernel density estimator and will be left to the reader.

# D  Proof of theorems on bias and variance

**Lemma D.1.** *Assume that $U(x, y)$ is any arbitrary functional which satisfies*

$$(i) \sup_{x \in (\epsilon_0, \epsilon_1)} |U(x, y)| = G_0 < \infty,$$

$$(ii) \sup_{x \in (q_l, q_u)} |U(x, y)|\mathcal{C}(k) = G_1 < \infty,$$

$$(iii) \mathbb{E}[\sup_{x \in (p_l, p_u)} |U(x/\mathbf{p}, y)|] = G_2 < \infty.$$

*Let $\mathbf{Z}$ denote $\mathbf{X}_i$ for some fixed $i \in \{1, .., N\}$. Let $\zeta_{\mathbf{Z}}$ be any random variable which almost surely lies in the range $(f(\mathbf{Z}), \tilde{\mathbf{f}}_k(\mathbf{Z}))$. Then,*

$$\mathbb{E}[|U(\zeta_{\mathbf{Z}}, \mathbf{Z})|] < \infty.$$

*Proof.* We will show that the conditional expectation $\mathbb{E}[|U(\zeta_Z, Z)| \mid \mathcal{X}_N] < \infty$. Because $0 < \epsilon_0 < f(X) < \epsilon_\infty < \infty$ by $(\mathcal{A}.1)$, it immediately follows that

$$\mathbb{E}[|U(\zeta_{\mathbf{Z}}, \mathbf{Z})|] = \mathbb{E}[\mathbb{E}[|U(\zeta_Z, Z)| \mid \mathcal{X}_N]] < \infty.$$

For fixed $\mathcal{X}_N$, $Z \in \mathcal{I}_N$ or $Z \in \mathcal{B}_N$. These two cases are handled seperately.

**Case 1:** $Z \in \mathcal{I}_N$  In this case, $\tilde{\mathbf{f}}_k(Z) = \hat{\mathbf{f}}_k(Z)$. By (76) and $(\mathcal{A}.1)$, we know that if $\natural(Z)$ holds, $p_l/\mathbf{P}(Z) < \hat{\mathbf{f}}_k(Z) < p_u/\mathbf{P}(Z)$. On the other hand, if $\natural^c(Z)$ holds, by (78) and $(\mathcal{A}.1)$,

$q_l < \hat{\mathbf{f}}_k(Z) < q_u$. This therefore implies that if $\natural(Z)$ holds, $\min\{\epsilon_0, p_l/\mathbf{P}(Z)\} < \zeta_Z < \max\{\epsilon_\infty, p_u/\mathbf{P}(Z)\}$ and if $\natural^c(Z)$ holds, $\min\{\epsilon_0, q_l\} < \zeta_Z < \max\{\epsilon_\infty, q_u\}$. Then,

$$
\begin{aligned}
\mathbb{E}[|U(\zeta_Z, Z)| \mid \mathcal{X}_N] &= \mathbb{E}[1_{\natural(Z)}|U(\zeta_Z, Z)| \mid \mathcal{X}_N] + \mathbb{E}[1_{\natural^c(Z)}|U(\zeta_Z, Z)| \mid \mathcal{X}_N] \\
&\leq G_0 + \mathbb{E}[1_{\natural(Z)} \sup_{x \in (p_l, p_u)} |U(x/\mathbf{P}(Z), Z)|] + \max\{G_0, G_1/\mathcal{C}(k)\}(1 - Pr(\natural(Z))) \\
&\leq G_0 + \mathbb{E}[\sup_{x \in (p_l, p_u)} |U(x/\mathbf{P}(Z), Z)|] + \max\{G_0, G_1/\mathcal{C}(k)\}(1 - Pr(\natural(Z))) \\
&= G_0 + G_2 + \max\{G_1/\mathcal{C}(M), G_0\}\mathcal{C}(k) \\
&= G_0 + G_2 + \max\{G_1, G_0\mathcal{C}(k)\} < \infty
\end{aligned}
\tag{132}
$$

where the final step follows from the fact that $\mathcal{C}(k) = o(1)$.

**Case 2:** $Z \in \mathcal{B}_N$  If $Z \in \mathcal{B}_N$, let $Y_n$ be the nearest neighbor of $Z$ in the set $\mathcal{I}_N$. Then,

$$
\tilde{\mathbf{f}}_k(Z) = \hat{\mathbf{f}}_k(Y_n)
\tag{133}
$$

This implies that we can now condition on the event $\natural(Y_n)$, and follow the exact procedure as in case 1 to obtain

$$
\begin{aligned}
\mathbb{E}[|U(\zeta_Z, Z)| \mid \mathcal{X}_N] &= \mathbb{E}[1_{\natural(Y_n)}|U(\zeta_Z, Z)| \mid \mathcal{X}_N] + \mathbb{E}[1_{\natural^c(Y_n)}|U(1/\zeta_Z, Z)| \mid \mathcal{X}_N] \\
&\leq G_0 + G_2 + \max\{G_1, G_0\mathcal{C}(k)\} < \infty
\end{aligned}
\tag{134}
$$

where the final step follows from the fact that $\mathcal{C}(k) = o(1)$. This concludes the proof.

$\square$

**Proof of Theorem 3.1**.

*Proof.* Using the continuity of $g'''(x, y)$, construct the following third order Taylor series of $g(\tilde{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z})$ around the conditional expected value $\mathbb{E}_Z[\tilde{\mathbf{f}}_k(\mathbf{Z})] = \mathbb{E}[\tilde{\mathbf{f}}_k(\mathbf{Z}) \mid \mathbf{Z}]$.

$$
\begin{aligned}
g(\tilde{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) = g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z}) + g'(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})\mathbf{e}(\mathbf{Z}) \\
+ \frac{1}{2}g''(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})\mathbf{e}^2(\mathbf{Z}) + \frac{1}{6}g^{(3)}(\zeta_{\mathbf{Z}}, \mathbf{Z})\mathbf{e}^3(\mathbf{Z}),
\end{aligned}
$$

where $\zeta_{\mathbf{Z}} \in (\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \tilde{\mathbf{f}}_k(\mathbf{Z}))$ is defined by the mean value theorem. This gives

$$
\begin{aligned}
&\mathbb{E}[(g(\tilde{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z}))] \\
&= \mathbb{E}\left[\frac{1}{2}g''(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})\mathbf{e}^2(\mathbf{Z})\right] + \mathbb{E}\left[\frac{1}{6}g^{(3)}(\zeta_{\mathbf{Z}}, \mathbf{Z})\mathbf{e}^3(\mathbf{Z})\right]
\end{aligned}
$$

Let $\Delta(\mathbf{Z}) = \frac{1}{6}g^{(3)}(\zeta_{\mathbf{Z}}, \mathbf{Z})$. Direct application of Lemma D.1 in conjunction with assumptions $(\mathcal{A}.5)$, $(\mathcal{A}.6)$ implies that $\mathbb{E}[\Delta^2(\mathbf{Z})] = O(1)$. By Cauchy-Schwarz and assumption $(\mathcal{A}.4)$ applied to (125) for the choice $q = 6$,

$$
\left|\mathbb{E}\left[\frac{1}{6}\Delta(\mathbf{Z})\mathbf{e}^3(\mathbf{Z})\right]\right| \leq \sqrt{\mathbb{E}\left[\frac{1}{36}\Delta^2(\mathbf{Z})\right]\mathbb{E}[\mathbf{e}^6(\mathbf{Z})]} = o\left(\frac{1}{k}\right) + O(N\mathcal{C}(k)).
$$

77

By observing that the density estimates $\{\tilde{\mathbf{f}}_k(\mathbf{X}_i)\}, i = 1, \ldots, N$ are identical, we therefore have

$$\mathbb{E}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] - G(f) = \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})]$$
$$= \mathbb{E}[g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})] + \mathbb{E}\left[\frac{1}{2}g''(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})\mathbf{e}^2(\mathbf{Z})\right] + o(1/k) + O(N\mathcal{C}(k)).$$

By (131) and (125) for the choice $q = 2$, in conjunction with assumption $(\mathcal{A}.4)$, this implies that

$$
\begin{aligned}
\mathbb{E}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] - G(f) &= \mathbb{E}[g'(f(\mathbf{Z}), \mathbf{Z})h(\mathbf{Z})]\left(\frac{k}{M}\right)^{2/d} + \mathbb{E}[1_{\{\mathbf{Z}\in\mathcal{S}-\mathcal{S}_I\}}(g(f(\mathbf{Z}_{-1}), \mathbf{Z}_{-1}) - g(f(\mathbf{Z}), \mathbf{Z}))] \\
&\quad + \mathbb{E}[f^2(\mathbf{Z})g''(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})/2]\left(\frac{1}{k}\right) + O(N\mathcal{C}(k)) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right) \\
&= \mathbb{E}[g'(f(\mathbf{Z}), \mathbf{Z})h(\mathbf{Z})]\left(\frac{k}{M}\right)^{2/d} + \mathbb{E}[1_{\{\mathbf{Z}\in\mathcal{S}-\mathcal{S}_I\}}(g(f(\mathbf{Z}_{-1}), \mathbf{Z}_{-1}) - g(f(\mathbf{Z}), \mathbf{Z}))] \\
&\quad + \mathbb{E}[f^2(\mathbf{Z})g''(f(\mathbf{Z}), \mathbf{Z})/2]\left(\frac{1}{k}\right) + O(N\mathcal{C}(k)) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right) \\
&= c_1\left(\frac{k}{M}\right)^{2/d} + c_2\left(\frac{1}{k}\right) + c_3 + O(N\mathcal{C}(k)) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right),
\end{aligned}
$$

where the last but one step follows because, by (81) and (124), we know $\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})] = f(\mathbf{Z}) + o(1)$. This in turn implies $\mathbb{E}[f^2(\mathbf{Z})g''(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})/2] = \mathbb{E}[f^2(\mathbf{Y})g''(f(\mathbf{Y}), \mathbf{Y})/2]$. Finally, by assumption $(\mathcal{A}.5)$ and $(\mathcal{A}.2)$, the leading constants $c_1$ and $c_2$ are bounded. We have also shown in equation (130) that $c_3 = O((k/M)^{2/d})$. This concludes the proof.

$\square$

**Proof of Theorem 5.1**

*Proof.* Let $\mathbf{X}$ denote $\mathbf{X}_i$ for some fixed $i \in \{1, .., N\}$. Also, let $\mathbf{X}_{-1} = \arg\min_{x\in\mathcal{S}_I} d(x, \mathbf{X})$.

Using (82), we can derive the following in an identical manner to (131):

$$
\begin{aligned}
\mathbb{B}(\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)) &= \mathbb{E}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)] - \int g(f(x),x)f(x)dx \\
&= (\mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{Z}),\mathbf{Z})] - g_2(k,M))/g_1(k,M) - \int g(f(x),x)f(x)dx \\
&= \mathbb{E}[\mathbb{E}[(g(\tilde{\mathbf{f}}_k(\mathbf{Z}),\mathbf{X}) - g_2(k,M))/g_1(k,M) \mid \mathcal{X}_N]] - \int g(f(x),x)f(x)dx \\
&= \mathbb{E}[\mathbb{E}[(g(\tilde{\mathbf{f}}_k(\mathbf{X}),\mathbf{X}) - g_2(k,M))/g_1(k,M) \mid \mathcal{X}_N], X \in \mathcal{I}_N] \\
&\quad +\mathbb{E}[\mathbb{E}[(g(\tilde{\mathbf{f}}_k(\mathbf{X}),\mathbf{X}) - g_2(k,M))/g_1(k,M) \mid \mathcal{X}_N], X \in \mathcal{B}_N] \\
&\quad -\int g(f(x),x)f(x)dx \\
&= \mathbb{E}[g(f(\mathbf{X}),\mathbf{X}) + \frac{g'(f(\mathbf{X}),\mathbf{X})h(\mathbf{X})}{g_1(k,M)}(k/M)^{2/d} \\
&\quad +\frac{1_{\{\mathbf{X}\in\mathcal{S}-\mathcal{S}'\}}}{g_1(k,M)}(g(f(\mathbf{X}_{-1}),\mathbf{X}_{-1}) - g(f(\mathbf{X}),\mathbf{X})) \\
&\quad +o((k/M)^{2/d}) + O(N\mathcal{C}(k))] - \int g(f(x),x)f(x)dx \\
&= \frac{c_1}{g_1(k,M)}\left(\frac{k}{M}\right)^{2/d} + \frac{c_3}{g_1(k,M)} + o\left(\left(\frac{k}{M}\right)^{2/d}\right) + O(N\mathcal{C}(k)).
\end{aligned}
$$

Because we assume the logarithmic growth condition $k = O((\log(M))^{2/(1-\delta)})$, it follows that $O(N\mathcal{C}(k)) = O(N/M^3) = o(1/T)$. Also, by (8), $g_1(k,M) = 1 + o(1)$. This implies that

$$
\mathbb{B}(\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)) = c_1\left(\frac{k}{M}\right)^{2/d} + c_3 + o\left(\left(\frac{k}{M}\right)^{2/d}\right). \tag{135}
$$

$\square$

**Proof of Theorem 3.2 and Theorem 5.2.**

*Proof.* By the continuity of $g^{(\lambda)}(x,y)$, we can construct the following Taylor series of $g(\tilde{\mathbf{f}}_k(\mathbf{Z}),\mathbf{Z})$ around the conditional expected value $\mathbb{E}_Z[\tilde{\mathbf{f}}_k(\mathbf{Z})]$.

$$
\begin{aligned}
g(\tilde{\mathbf{f}}_k(\mathbf{Z}),\mathbf{Z}) &= g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})],\mathbf{Z}) + g'(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})],\mathbf{Z})\mathbf{e}(\mathbf{Z}) \\
&\quad + \left(\sum_{i=2}^{\lambda-1} \frac{g^{(i)}(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})],\mathbf{Z})}{i!}\mathbf{e}^i(\mathbf{Z})\right) + \frac{g^{(\lambda)}(\xi_{\mathbf{Z}},\mathbf{Z})}{\lambda!}\mathbf{e}^\lambda(\mathbf{Z}),
\end{aligned}
$$

where $\xi_{\mathbf{Z}} \in (g(\mathbb{E}_Z[\tilde{\mathbf{f}}_k(\mathbf{Z})],g(\tilde{\mathbf{f}}_k(\mathbf{Z})))$. Denote $(g^\lambda(\xi_{\mathbf{Z}},\mathbf{Z}))/\lambda!$ by $\Psi(\mathbf{Z})$. Further define the

operator $\mathcal{M}(\mathbf{Z}) = \mathbf{Z} - \mathbb{E}[\mathbf{Z}]$ and

$$
\begin{aligned}
p_i &= \mathcal{M}(g(\mathbb{E}_{\mathbf{X}_i}[\tilde{\mathbf{f}}_k(\mathbf{X}_i)], \mathbf{X_i})), \\
q_i &= \mathcal{M}(g'(\mathbb{E}_{\mathbf{X}_i}[\tilde{\mathbf{f}}_k(\mathbf{X}_i)], \mathbf{X_i})\mathbf{e}(\mathbf{X_i})), \\
r_i &= \mathcal{M}\left( \sum_{i=2}^{\lambda} \frac{g^{(i)}(\mathbb{E}_{\mathbf{X}_i}[\tilde{\mathbf{f}}_k(\mathbf{X}_i)], \mathbf{X_i})}{i!} \mathbf{e}^i(\mathbf{X_i}) \right) \\
s_i &= \mathcal{M}\left( \Psi(\mathbf{X_i})\mathbf{e}^{\lambda}(\mathbf{X_i}) \right)
\end{aligned}
$$

The variance of the estimator $\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)$ is given by

$$
\begin{aligned}
\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] &= \mathbb{E}[(\hat{\mathbf{G}}(f) - \mathbb{E}[\hat{\mathbf{G}}(f)])^2] \\
&= \frac{1}{N}\mathbb{E}\left[(p_1 + q_1 + r_1 + s_1)^2\right] \\
&+ \frac{N-1}{N}\mathbb{E}[(p_1 + q_1 + r_1 + s_1)(p_2 + q_2 + r_2 + s_2)].
\end{aligned}
$$

Because $\mathbf{X}_1$, $\mathbf{X}_2$ are independent, we have $\mathbb{E}[(p_1)(p_2 + q_2 + r_2 + s_2)] = 0$. Furthermore,

$$
\mathbb{E}\left[(p_1 + q_1 + r_1 + s_1)^2\right] = \mathbb{E}[p_1{}^2] + o(1) = \mathbb{V}[g(\mathbb{E}_{\mathbf{Z}}[\hat{\mathbf{f}}(\mathbf{Z})], \mathbf{Z})] + o(1).
$$

From assumption $(\mathcal{A}.4)$ applied to (125) and (126), in conjunction with assumption $(\mathcal{A}.3)$, it follows that

- $\mathbb{E}[p_1{}^2] = \mathbb{V}[g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})]$

- $\mathbb{E}[q_1 q_2] = \mathbb{V}[g'(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})f(\mathbf{Z})]\left(\frac{1}{M}\right) + o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k))$

- $\mathbb{E}[q_1 r_2] = \sum_{i=2}^{\lambda-1} O\left(\frac{1}{k^{((1+i)\delta/2-1)}M}\right) + O\left(\frac{\lambda(k_M^{2/d}+1/M)}{M}\right) + O(N\mathcal{C}(k)) = o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k))$

- $\mathbb{E}[r_1 r_2] = \sum_{i_1=2}^{\lambda-1} \sum_{i_2=2}^{\lambda-1} O\left(\frac{1}{k^{((i_1+i_2)\delta/2-1)}M}\right) + O\left(\frac{\lambda^2(k_M^{2/d}+1/M)}{M}\right) + O(N\mathcal{C}(k)) = o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k))$

Since $q_1$ and $s_2$ are 0 mean random variables

$$
\begin{aligned}
\mathbb{E}[q_1 s_2] &= \mathbb{E}\left[q_1\Psi(\mathbf{X_2})(\hat{\mathbf{f}}(\mathbf{X_2}) - \mathbb{E}_{\mathbf{X}_2}[\tilde{\mathbf{f}}_k(\mathbf{X_2})])^{\lambda}\right] \\
&= \mathbb{E}\left[q_1\Psi(\mathbf{X_2})(\hat{\mathbf{f}}(\mathbf{X_2}) - \mathbb{E}_{\mathbf{X}_2}[\tilde{\mathbf{f}}_k(\mathbf{X_2})])^{\lambda}\right] \\
&\leq \sqrt{\mathbb{E}\left[\Psi^2(\mathbf{X_2})\right]\mathbb{E}\left[q_1^2(\hat{\mathbf{f}}(\mathbf{X_2}) - \mathbb{E}_{\mathbf{X}_2}[\tilde{\mathbf{f}}_k(\mathbf{X_2})])^{2\lambda}\right]} \\
&= \sqrt{\mathbb{E}\left[\Psi^2(\mathbf{Z})\right]}\left(o\left(\frac{1}{k^{\lambda}}\right) + O(N\mathcal{C}(k))\right)
\end{aligned}
$$

Direct application of Lemma D.1 in conjunction with assumptions $(\mathcal{A}.5)$, $(\mathcal{A}.6)$ implies that $\mathbb{E}\left[\Psi^2(\mathbf{Z})\right] = O(1)$. Note that from assumption $(\mathcal{A}.3)$, $o\left(\frac{1}{k^\lambda}\right) = o(1/M)$ . In a similar manner, it can be shown that $\mathbb{E}[r_1 s_2] = o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k))$ and $\mathbb{E}[s_1 s_2] = o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k))$. Finally, by (81) and (124), we know $\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})] = \mathbb{E}[\tilde{\mathbf{f}}_k(\mathbf{Z})] = f(\mathbf{Z}) + o(1)$. This implies that

$$
\begin{aligned}
\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)] &= \frac{1}{N}\mathbb{E}\left[p_1{}^2\right] + \frac{(N-1)}{N}\mathbb{E}[q_1 q_2] + O(N\mathcal{C}(k)) + o\left(\frac{1}{M} + \frac{1}{N}\right) \\
&= \mathbb{V}[g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})]\left(\frac{1}{N}\right) + \mathbb{V}[g'(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})f(\mathbf{Z})]\left(\frac{1}{M}\right) + O(N\mathcal{C}(k)) + o\left(\frac{1}{M} + \frac{1}{N}\right) \\
&= \mathbb{V}[g(f(\mathbf{Z}), \mathbf{Z})]\left(\frac{1}{N}\right) + \mathbb{V}[g'(f(\mathbf{Z}), \mathbf{Z})f(\mathbf{Z})]\left(\frac{1}{M}\right) + O(N\mathcal{C}(k)) + o\left(\frac{1}{M} + \frac{1}{N}\right) \\
&= c_4\left(\frac{1}{N}\right) + c_5\left(\frac{1}{M}\right) + O(N\mathcal{C}(k)) + o\left(\frac{1}{M} + \frac{1}{N}\right),
\end{aligned}
$$

where the last but one step follows because, by (81) and (124), we know $\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})] = f(\mathbf{Z}) + o(1)$. This in turn implies $\mathbb{V}[g(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})] = \mathbb{V}[g(f(\mathbf{Z}), \mathbf{Z})]$ and $\mathbb{V}[g'(\mathbb{E}_{\mathbf{Z}}[\tilde{\mathbf{f}}_k(\mathbf{Z})], \mathbf{Z})f(\mathbf{Z})] = \mathbb{V}[g'(f(\mathbf{Z}), \mathbf{Z})f(\mathbf{Z})]$. Finally, by assumptions $(\mathcal{A}.5)$ and $(\mathcal{A}.2)$, the leading constants $c_4$ and $c_5$ are bounded. This concludes the proof of Theorem 3.2.

Under the logarithmic growth condition $k = O((\log(M))^{2/(1-\delta)})$, $g_2(k, M) = o(1)$ and $g_1(k, M) = 1 + o(1)$ by assumption (8). Theorem 5.2 follows by observing that $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) = (\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - g_1(k, M))/g_2(k, M)$ □

**Bias of Baryshnikov's estimator: Proof of equation (5)**

*Proof.* We will first prove that

$$\mathbb{B}(\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_k)) = \Theta((k/M)^{1/d} + 1/k), \tag{136}$$

Because the standard $k$-NN density estimate $\hat{\mathbf{f}}_{kS}(\mathbf{X}_i)$ is identical to the partitioned $k$-NN density estimate $\hat{\mathbf{f}}_k(\mathbf{X}_i)$ defined on the partition $\{\mathbf{X}_i\}$ and $\{\mathbf{X}_1, .., \mathbf{X}_T\} - \{\mathbf{X}_i\}$, it follows that

$$\mathbb{B}(\tilde{\mathbf{G}}_N(\hat{\mathbf{f}}_{kS})) = \Theta((k/T)^{1/d} + 1/k). \tag{137}$$

From the definition of set $\mathcal{S}'$ in section B.2.1, we can choose the set $\mathcal{S}'$, such that $Pr(\mathbf{Z} \notin \mathcal{S}') = O((k/M)^{1/d})$.

$$
\begin{aligned}
\mathbb{E}[\hat{\mathbf{G}}_N(\hat{\mathbf{f}}_k)] - G(f) &= \mathbb{E}[g(\hat{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})] \\
&= \mathbb{E}[1_{\{\mathbf{Z} \in \mathcal{S}'\}} g(\hat{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})] + \mathbb{E}[1_{\{\mathbf{Z} \in \mathcal{S} - \mathcal{S}'\}} g(\hat{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})] \\
&= I + II
\end{aligned} \tag{138}
$$

Using the exact same method as in the Proof of Theorem 3.1, using (81) and (116), and the fact that $Pr(\mathbf{Z} \notin \mathcal{S}') = O((k/M)^{1/d}) = o(1)$, we have

$$I = \mathbb{E}[g'(f(\mathbf{Z}), \mathbf{Z})h(\mathbf{Z})]\left(\frac{k}{M}\right)^{2/d} + \mathbb{E}[f^2(\mathbf{Z})g''(f(\mathbf{Z}), \mathbf{Z})/2]\left(\frac{1}{k}\right) + O(\mathcal{C}(k)) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right),$$

Because we assume that $g$ satisfies assumption $(\mathcal{A}.6)$, from the proof of Lemma D.1, for $Z \in \mathcal{S} - \mathcal{S}'$, we have $\mathbb{E}[g(\hat{\mathbf{f}}_k(Z), Z) - g(f(Z), Z)] = O(1)$. This implies that,

$$
\begin{aligned}
II &= \mathbb{E}[1_{\{\mathbf{Z} \in \mathcal{S} - \mathcal{S}'\}} g(\hat{\mathbf{f}}_k(\mathbf{Z}), \mathbf{Z}) - g(f(\mathbf{Z}), \mathbf{Z})] \\
&= \mathbb{E}\left[\mathbb{E}[g(\hat{\mathbf{f}}_k(Z), Z) - g(f(Z), Z)] \mid 1_{\{\mathbf{Z} \in \mathcal{S} - \mathcal{S}'\}}\right] \times Pr(\mathbf{Z} \notin \mathcal{S}') \\
&= O(1) \times O((k/M)^{1/d}) = O((k/M)^{1/d}).
\end{aligned}
\tag{139}
$$

This concludes the proof.

$\square$

# E    Asymptotic normality

Define the random variables $\{\mathbf{Y}_{M,i}; i = 1, \ldots, N\}$ for any fixed $M$

$$
\mathbf{Y}_{M,i} = \frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)]}{\sqrt{\mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)]}},
$$

and define the sum $\mathbf{S}_{\mathbf{N},\mathbf{M}}$

$$
\mathbf{S}_{N,M} = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} \mathbf{Y}_{M,i},
$$

where the indices $N$ and $M$ explicitly stress the dependence of the sum $\mathbf{S}_{N,M}$ on the number of random variables $N + M$. Observe that the random variables $\{\mathbf{Y}_{M,i}; i = 1, \ldots, N\}$ belong to an 0 mean, unit variance, interchangeable process [5] for all values of $M$. To establish the CLT for $\mathbf{S}_{N,M}$, we will exploit the fact the random variables $\{\mathbf{Y}_{M,i}; i = 1, \ldots, N\}$ are interchangeable by appealing to DeFinetti's theorem, which we describe below.

## E.1    De Finetti's Theorem

Let $\mathcal{F}$ be the class of one dimensional distribution functions and for each pair of real numbers $x$ and $y$ define $\mathcal{F}(x, y) = \{F \in \mathcal{F}|F(x) \le y\}$. Let $\mathcal{B}$ be the Borel field of subsets of $\mathcal{F}$ generated by the class of sets $\mathcal{F}(x, y)$. Then De Finetti's theorem asserts that for any interchangeable process $\{\mathbf{Z}_i\}$ there exists a probability measure $\mu$ defined on $\mathcal{B}$ such that

$$
Pr\{\mathbf{B}\} = \int_{\mathcal{F}} Pr_F\{\mathbf{B}\} d\mu(F),
\tag{140}
$$

for any Borel measurable set defined on the sample space of the sequence $\{\mathbf{Z}_i\}$. Here $Pr\{\mathbf{B}\}$ is the probability of the event $\mathbf{B}$ and $Pr_F\{\mathbf{B}\}$ is the probability of the event $B$ under the assumption that component random variables $\mathbf{X}_i$ of the interchangeable process are independent and identically distributed with distribution $F$.

## E.2 Necessary and Sufficient conditions for CLT

For each $F \in \mathcal{F}$ define $m(F)$ and $\sigma^2(F)$ as $m(F) = \int_{-\infty}^{\infty} x dF(x)$, $\sigma(F) = \int_{-\infty}^{\infty} x^2 dF(x) - 1$ and for all real numbers $m$ and non-negative real numbers $\sigma^2$ let $\mathcal{F}_{m,\sigma^2}$ be the set of $F \in \mathcal{F}$ for which $m(F) = m$ and $\sigma^2(F) = \sigma^2$.

Let $\{\mathbf{Z}_i; i = 1, 2, \ldots\}$ be an interchangeable stochastic process with 0 mean and variance 1. Blum *etal* [5] showed that the random variable $\mathbf{S}_N = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} \mathbf{Z}_i$ converges in distribution to $N(0,1)$ if and only if $\mu(\mathcal{F}_{0,0}) = 1$. Furthermore, they show that the condition $\mu(\mathcal{F}_{0,0}) = 1$ is equivalent to the condition that $Cov(\mathbf{Z}_1, \mathbf{Z}_2) = 0$ and $Cov(\mathbf{Z}_1^2, \mathbf{Z}_2^2) = 0$. We will *extend* Blum *etal*'s results to interchangeable processes where $Cov(\mathbf{Z}_1, \mathbf{Z}_2) = o(1)$ and $Cov(\mathbf{Z}_1^2, \mathbf{Z}_2^2) = o(1)$.

In particular, we will show that $Cov(\mathbf{Y}_{M,1}, \mathbf{Y}_{M,2})$ and $Cov(\mathbf{Y}_{M,1}^2, \mathbf{Y}_{M,2}^2)$ are $O(1/M)$. Subsequently we will show that the random variable $\mathbf{S}_{N,M} = \frac{1}{\sqrt{N}} \sum_{i=1}^{N} \mathbf{Y}_{M,i}$ converges in distribution to $N(0,1)$ and conclude that Theorem 3.3 holds.

## E.3 CLT for Asymptotically Uncorrelated processes

Let $\mathbf{X}$ be a random variable with density $f$. In the proof of Theorem 3.2, we showed that

$$
\begin{aligned}
Cov(\mathbf{Y}_{M,i}, \mathbf{Y}_{M,j}) &= \frac{Cov(g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i), g(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j))}{\sqrt{\mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)] \mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j)]}} \\
&= \frac{Cov(p_i + q_i + r_i + s_i, p_j + q_j + r_j + s_j)}{\sqrt{\mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)] \mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j)]}} \\
&= \frac{Cov(p_i + q_i + r_i + s_i, p_j + q_j + r_j + s_j)}{\sqrt{\mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)] \mathbb{V}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j)]}} \\
&= \frac{\mathbb{V}(g'(f(\mathbf{X}), \mathbf{X}) f(\mathbf{X}))}{\mathbb{V}[g(f(\mathbf{X}_i), \mathbf{X}_i)]} \left(\frac{1}{M}\right) + o\left(\frac{1}{M}\right) + O(N\mathcal{C}(k)) \\
&= \frac{\mathbb{V}(g'(f(\mathbf{X}), \mathbf{X}) f(\mathbf{X}))}{\mathbb{V}[g(f(\mathbf{X}_i), \mathbf{X}_i)]} \left(\frac{1}{M}\right) + o\left(\frac{1}{M}\right),
\end{aligned}
\tag{141}
$$

where the last but one step follows by observing that $N\mathcal{C}(k)/M \to 0$ under the logarithmic growth condition $k = O((\log(M))^{2/(1-\delta)})$. Define the function $d(x, y) = g(x, y)(g(x, y) - c)$, where the constant $c = \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}), \mathbf{X})]$. Then, similar to the derivation of (141), we have,

$$
\begin{aligned}
Cov(\mathbf{Y}_{M,i}^2, \mathbf{Y}_{M,j}^2) &= \frac{Cov(d(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i), d(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j))}{\sqrt{\mathbb{V}[d(\tilde{\mathbf{f}}_k(\mathbf{X}_i), \mathbf{X}_i)] \mathbb{V}[d(\tilde{\mathbf{f}}_k(\mathbf{X}_j), \mathbf{X}_j)]}} \\
&= \frac{\mathbb{V}(d'(f(\mathbf{X}), \mathbf{X}) f(\mathbf{X}))}{\mathbb{V}[d(f(\mathbf{X}_i), \mathbf{X}_i)]} \left(\frac{1}{M}\right) + o\left(\frac{1}{M}\right).
\end{aligned}
\tag{142}
$$

83

**Proof of Theorem 3.3 and Theorem 5.3**.

*Proof.* Let $\delta_\mu(M)$ and $\delta_\sigma(M)$ be a strictly positive functions parameterized by $M$ such that $\delta_\mu(M) = o(1); \frac{1}{M\delta_\mu(M)} = o(1)$, $\delta_\sigma(M) = o(1); \frac{1}{M\delta_\sigma(M)} = o(1)$. Denote the set of $F \in \mathcal{F}$ with $\mathcal{F}_{m,\delta,M} := \{m^2(F) \geq \delta_\mu(M)\}$; $\mathcal{F}_{\sigma,\delta,M} := \{\sigma^2(F) \geq \delta_\sigma(M)\}$; $\mathcal{F}^*_{m,\delta,M} := \{m^2(F) \in (0, \delta_\mu(M))\}$ and $\mathcal{F}^*_{\sigma,\delta,M} := \{\sigma^2(F) \in (0, \delta_\sigma(M))\}$. Denote the measures of these sets by $\mu_{m,\delta,M}$, $\mu_{\sigma,\delta,M}$, $\mu^*_{m,\delta,M}$ and $\mu^*_{\sigma,\delta,M}$ respectively. We have from (140) that

$$
\begin{aligned}
\int_{\mathcal{F}} m^2(F) d\mu(F) &= Cov(\mathbf{Y}_{M,i}, \mathbf{Y}_{M,j}) \\
\int_{\mathcal{F}} \sigma^2(F) d\mu(F) &= \int_{\mathcal{F}} [\mathbb{E}_F[\mathbf{Z}^2 - 1]]^2 d\mu(F) = Cov(\mathbf{Y}^2_{M,i}, \mathbf{Y}^2_{M,j}).
\end{aligned}
\tag{143}
$$

Applying the Chebyshev inequality, we get

$$
\begin{aligned}
\delta_\mu(M)\mu_{m,\delta,M} &\leq Cov(\mathbf{Y}_{M,i}, \mathbf{Y}_{M,j}), \\
\delta_\sigma(M)\mu_{\sigma,\delta,M} &\leq Cov(\mathbf{Y}^2_{M,i}, \mathbf{Y}^2_{M,j}).
\end{aligned}
$$

Because the covariances decay at $O(1/M)$, $\mu_{m,\delta,M}$ and $\mu_{\sigma,\delta,M} \to 0$ as $M \to \infty$. From the definition of $\mathcal{F}^*_{m,\delta,M}$ and $\mathcal{F}^*_{\sigma,\delta,M}$, we also have that $\mu^*_{m,\delta,M}$ and $\mu^*_{\sigma,\delta,M} \to 0$ as $M \to \infty$. We also have

$$
1 - (\mu_{m,\delta,M} + \mu_{\sigma,\delta,M} + \mu^*_{m,\delta,M} + \mu^*_{\sigma,\delta,M}) \leq \mu(\mathcal{F}_{0,0}) \leq 1,
$$

and therefore

$$
\lim_{M \to \infty} \mu(\mathcal{F}_{0,0}) = 1.
\tag{144}
$$

We will now show that $\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) = (\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - \mathbb{E}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)])/(\sqrt{\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)]})$ converges weakly to

$\mathbb{N}(0,1)$. Denote $g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i)$ by $\mathbf{g}_i$. Observe that

$$\lim_{\Delta\to 0} Pr\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\} = \lim_{\Delta\to 0}\int_{\mathcal{F}} Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}d\mu(F)$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}d\mu(F) + \lim_{\Delta\to 0}\int_{\mathcal{F}} 1_{\{F\in\mathcal{F}-\mathcal{F}_{0,0}\}}Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}d\mu(F)$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}d\mu(F) + \int_{\mathcal{F}} \lim_{\Delta\to 0}\left(1_{\{F\in\mathcal{F}-\mathcal{F}_{0,0}\}}Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}\right)d\mu(F) \tag{145}$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\{\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)\le\alpha\}d\mu(F) \tag{146}$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\left\{\frac{1}{N}\sum_{i=1}^{N}\left(\frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i)]}{\sqrt{\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)]}}\right)\le\alpha\right\}d\mu(F)$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\left\{\frac{1}{N}\sum_{i=1}^{N}\left(\frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i)]}{\sqrt{\mathbb{V}[\mathbf{g}_i]/N + ((N-1)/N)Cov[\mathbf{g}_i,\mathbf{g}_j]}}\right)\le\alpha\right\}\int_{\mathcal{F}_{0,0}}d\mu(F)$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\left\{\frac{1}{N}\sum_{i=1}^{N}\left(\frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i)]}{\sqrt{\mathbb{V}[\mathbf{g}_i]/N + ((N-1)/N)\sqrt{\mathbb{V}[\mathbf{g}_i]\mathbb{V}[\mathbf{g}_j]}Cov[\mathbf{Y}_{M,i},\mathbf{Y}_{M,j}]}}\right)\le\alpha\right\}d\mu(F)$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\left\{\frac{1}{N}\sum_{i=1}^{N}\left(\frac{g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i) - \mathbb{E}[g(\tilde{\mathbf{f}}_k(\mathbf{X}_i),\mathbf{X}_i)]}{\sqrt{\mathbb{V}[\mathbf{g}_i]/N}}\right)\le\alpha\right\}d\mu(F) \tag{147}$$

$$= \lim_{\Delta\to 0}\int_{\mathcal{F}_{0,0}} Pr_F\left\{\frac{1}{\sqrt{N}}\sum_{i=1}^{N}\mathbf{Y}_{M,i}\le\alpha\right\}d\mu(F)$$

$$= \int_{\mathcal{F}} \lim_{\Delta\to 0}\left(1_{\{F\in\mathcal{F}_{0,0}\}}Pr_F\left\{\frac{1}{\sqrt{N}}\sum_{i=1}^{N}\mathbf{Y}_{M,i}\le\alpha\right\}\right)d\mu(F)$$

$$= \int_{\mathcal{F}}\phi(\alpha)d\mu(F) = \phi(\alpha), \tag{148}$$

where $\phi(.)$ is the distribution function of a Gaussian random variable with mean 0 and variance 1. Step (D.6) follows from the Dominated Convergence theorem. By (144), $\lim_{\Delta\to 0} 1_{\{F\in\mathcal{F}-\mathcal{F}_{0,0}\}} = 0$ almost surely. This gives Step (D.7). Step (D.8) is obtained by observing that, by (143), $Cov[\mathbf{Y}_{M,i},\mathbf{Y}_{M,j}] = 0$ when $F \in \mathcal{F}_{0,0}$. The last step (D.9) follows from the CLT for sums of 0 mean, unit variance, i.i.d random variables and (144). This concludes the proof of Theorem 3.3.

To show Theorem 5.3, observe that under the logarithmic growth condition $k = O((\log(M))^{2/(1-\delta)})$, $g_2(k,M) = o(1)$ and $g_1(k,M) = 1 + o(1)$ by assumption (8). Since $\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) = (\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - g_1(k,M))/g_2(k,M)$, it follows that the asymptotic distribution of

$$\frac{\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k) - \mathbb{E}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)]}{\sqrt{\mathbb{V}[\hat{\mathbf{G}}_{N,BC}(\tilde{\mathbf{f}}_k)]}}$$

is equal to the asymptotic distribution of $\tilde{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) = (\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k) - \mathbb{E}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)])/(\sqrt{\mathbb{V}[\hat{\mathbf{G}}_N(\tilde{\mathbf{f}}_k)]})$.

$\square$

## E.4   Berry-Esseen bounds

We now establish Berry-Esseen bounds for the case where $\frac{N}{M} \to 0$. In particular, we assume that there exists a $\delta : 0 < \delta < 1$, such that $N = O(M^\delta)$. We also assume that the interchangeable process has finite absolute third order moment $E(|\mathbf{Z}_{M,i}|^3) = \rho_M < \infty \vee M$.

### E.4.1   Details

Define the subset $\tilde{F}$ of $F$ as follows: $\tilde{F} = F - \{F_{m,\delta,M} \bigcup F_{\sigma,\delta,M}\}$.

We recognize that for $F \in \tilde{F}$, we have

$$
-\sqrt{\delta_\mu(M)} \leq m(F) \leq \sqrt{\delta_\mu(M)},
$$
$$
-\sqrt{\delta_\sigma(M)} \leq \sigma(F) \leq \sqrt{\delta_\sigma(M)}.
$$

The mean and variance of $Y_{M,i}$ under the distribution $F$ are given by $m(F)$ and $\sigma(F) + \rho - m^2(F)$ respectively.

As in the previous section, let $\phi$ be the distribution function of a Gaussian random variable with 0 mean and $\rho$ variance.

**Lower bound**

$$
Pr\{\mathbf{S}_{N,M} \leq \alpha\} = \int_F Pr_F\{\mathbf{S}_{N,M} \leq \alpha\} d\mu(F)
$$
$$
\geq \int_{\tilde{F}} Pr_F\{\mathbf{S}_{N,M} \leq \alpha\} d\mu(F)
$$
$$
\geq \int_{\tilde{F}} \left[ \phi\left( \frac{\alpha - \sqrt{N}m(F)}{1 + (\sigma(F) - m^2(F))/\rho} \right) - \frac{C\kappa(F)}{(\sigma(F) + \rho - m^2(F))^3 \sqrt{N}} \right] d\mu(F)
$$
$$
\geq \phi\left( \frac{\alpha - \sqrt{N\delta_\mu(M)}}{1 + (\sqrt{\delta_\sigma(M)})/\rho} \right) \mu(\tilde{F}) - \int_{\tilde{F}} \frac{C\kappa(F)}{(\rho - \sqrt{\delta_\sigma(M)} - \delta_\mu(M))^3 \sqrt{N}} d\mu(F)
$$
$$
\geq \phi\left( \frac{\alpha - \sqrt{N\delta_\mu(M)}}{1 + (\sqrt{\delta_\sigma(M)})/\rho} \right) \mu(\tilde{F}) - \frac{C\kappa}{(\rho - \sqrt{\delta_\sigma(M)} - \delta_\mu(M))^3 \sqrt{N}}.
$$

**Upper bound**

Denote $\mu(\tilde{F}^c) := \tilde{\mu}$. We note that $\tilde{\mu} \leq \mu_{m,\delta,M} + \mu_{\sigma,\delta,M}$.

$$
\begin{aligned}
Pr\{\mathbf{S}_{N,M} \leq \alpha\} &= \int_F Pr_F\{\mathbf{S}_{N,M} \leq \alpha\} d\mu(F) \\
&\leq \int_{\tilde{F}} Pr_F\{\mathbf{S}_{N,M} \leq \alpha\} d\mu(F) + \tilde{\mu} \\
&\leq \int_{\tilde{F}} \left[ \phi\left( \frac{\alpha - \sqrt{N}m(F)}{1 + (\sigma(F) - m^2(F))/\rho} \right) + \frac{C\kappa(F)}{(\sigma(F) + \rho - m^2(F))^3 \sqrt{N}} \right] d\mu(F) + \tilde{\mu} \\
&\leq \phi\left( \frac{\alpha + \sqrt{N}\delta_\mu(M)}{1 - (\sqrt{\delta_\sigma(M)} + \delta_\mu(M))/\rho} \right) \mu(\tilde{F}) + \int_{\tilde{F}} \frac{C\kappa(F)}{(\rho + \sqrt{\delta_\sigma(M)})^3 \sqrt{N}} d\mu(F) + \tilde{\mu} \\
&\leq \phi\left( \frac{\alpha - \sqrt{N}\delta_\mu(M)}{1 - (\sqrt{\delta_\sigma(M)} + \delta_\mu(M))/\rho} \right) \mu(\tilde{F}) + \frac{C\kappa}{(\rho + \sqrt{\delta_\sigma(M)})^3 \sqrt{N}} + \mu_{m,\delta,M} + \mu_{\sigma,\delta,M} \\
&\leq \phi\left( \frac{\alpha - \sqrt{N}\delta_\mu(M)}{1 - (\sqrt{\delta_\sigma(M)} + \delta_\mu(M))/\rho} \right) \mu(\tilde{F}) + \frac{C\kappa}{(\rho + \sqrt{\delta_\sigma(M)})^3 \sqrt{N}} + \frac{1}{M\delta_\mu(M)} + \frac{1}{M\delta_\sigma(M)}.
\end{aligned}
$$

We have shown that the appropriately normalized sum $S_{N,M}$ converges in distribution to a normal random variable. Also for the case where $N$ grows slower than $M$, we have established Berry-Esseen type bounds on the error.

# F  Uniform kernel based plug-in estimator

In this section, we will state the main results concerning uniform kernel plug-in estimators. The proofs for these results rely on the properties of the uniform kernel density estimates established in Appendix A and proofs for equivalent results for the $k$-NN plug-in estimators. Let $\hat{\mathbf{f}}_u$ denote the boundary corrected uniform kernel density estimate. Denote the uniform kernel plug-in estimator by

$$
\hat{\mathbf{G}}_{\mathbf{u}}(f) = \left( \frac{1}{N} \sum_{i=1}^{N} g(\hat{\mathbf{f}}_u(\mathbf{X}_i), \mathbf{X}_i) \right). \tag{149}
$$

Let $\mathbf{Y}$ denote a random variable with density function $f$.

## F.1 Results

**Corollary F.1.** *Suppose that the density $f$, the functional $g$ and the density estimate $\hat{\mathbf{f}}_u$ satisfy the necessary conditions listed above. The bias of the plug-in estimator $\hat{\mathbf{G}}_u(f)$ is then given by*

$$B_u(f) \;=\; c_1 \left(\frac{k}{M}\right)^{2/d} + c_2 \left(\frac{1}{k}\right) + o\left(\frac{1}{k} + \left(\frac{k}{M}\right)^{2/d}\right),$$

*where $c_1 = \mathbb{E}[g'(f(\mathbf{Y}), \mathbf{Y})c(\mathbf{Y})]$, $c_2 = \mathbb{E}[g''(f(\mathbf{Y}), \mathbf{Y})f(\mathbf{Y})/2]$ are constants which depend on the underlying density $f$.*

**Corollary F.2.** *Suppose that the density $f$, the functional $g$ and the density estimate $\hat{\mathbf{f}}_u$ satisfy the necessary conditions listed above. The variance of the plug-in estimator $\hat{\mathbf{G}}_u(f)$ is given by*

$$\mathbb{V}_u(f) \;=\; c_4 \left(\frac{1}{N}\right) + c_5 \left(\frac{1}{M}\right) + o\left(\frac{1}{M} + \frac{1}{N}\right),$$

*where $c_4 = \mathbb{V}[g(f(\mathbf{Y}), \mathbf{Y})]$ and $c_5 = \mathbb{V}[f(\mathbf{Y})g'(f(\mathbf{Y}), \mathbf{Y})]$ are constants which depend on the underlying density $f$.*

**Corollary F.3.** *Suppose that the density $f$, the functional $g$ and the density estimate $\hat{\mathbf{f}}_u$ satisfy the necessary conditions listed above. Further suppose $\mathbb{E}[|g(f)|^3]$ is finite. The asymptotic distribution of the plug-in estimator $\hat{\mathbf{G}}_u(f)$ is given by*

$$\lim_{\Delta(k,N,M)\to 0} Pr\left(\frac{\hat{\mathbf{G}}_u(f) - \mathbb{E}[\hat{\mathbf{G}}_u(f)]}{\sqrt{\mathbb{V}[f(\mathbf{Y})g'(f(\mathbf{Y}), \mathbf{Y})]/N}} \leq \alpha\right) = Pr(\mathbf{Z} \leq \alpha),$$

*where $\mathbf{Z}$ is a standard normal random variable.*

# References

[1] I. Ahmad and Pi-Erh Lin. A nonparametric estimation of the entropy for absolutely continuous distributions (corresp.). *Information Theory, IEEE Transactions on*, 22(3):372 – 375, may 1976.

[2] Yu. Baryshnikov, M. D. Penrose, and J.E. Yukich. Gaussian limits for generalized spacings. *Ann. Appl. Probab.*, 19(1):158–185, 2009.

[3] P. J. Bickel and Y. Ritov. Estimating integrated squared density derivatives: Sharp best order of convergence estimates. *Sankhya: The Indian Journal of Statistics*, 50:381–393, October 1988.

[4] L. Birge and P. Massart. Estimation of integral functions of a density. *The Annals of Statistics*, 23(1):11–29, 1995.

[5] J.R. Blum, H. Chernoff, M. Rosenblatt, and H. Teicher. Central limit theorems for interchangeable processes. *Canadian Journal of Mathematics*, June 1957.

[6] Y. Chen, A. Wiesel, and A. O. Hero. Robust shrinkage estimation of high-dimensional covariance matrices. submitted to IEEE Trans. on Signal Process., preprint available in arXiv:1009.5331.

[7] R. C. H. Cheng and N. A. K. Amin. Estimating parameters in continuous univariate distributions with a shifted origin. *Journal of the Royal Statistical Society. Series B (Methodological)*, 11:394–403, 1983.

[8] C. I. Chow and C. N. Liu. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory*, 14:462–467, 1968.

[9] J.A. Costa, A. Girotra, and A.O. Hero. Estimating local intrinsic dimension with k-nearest neighbor graphs. In *2005 IEEE/SP 13th Workshop on Statistical Signal Processing*, pages 417–422, 2005.

[10] E. J. Dudewicz and E. C. van der Meulen. Entropy-based tests of uniformity. *Journal of the American Statistical Association*, 76:967–974, 1981.

[11] P. B. Eggermont and V. N. LaRiccia. Best asymptotic normality of the kernel density entropy estimator for smooth densities. *Information Theory, IEEE Transactions on*, 45(4):1321 –1326, May 1999.

[12] D. Evans. A law of large numbers for nearest neighbor statistics. *Proceedings of the Royal Society A*, 464:3175–3192, 2008.

[13] D. Evans, A. Jones, and W. M. Schmidt. Asymptotic moments of nearest neighbor distance distributions. *Proceedings of the Royal Society A*, 458:2839–2849, 2008.

[14] A.M. Farahmand, C. Sepesvari, and J-Y Audibert. Manifold-adaptive dimension estimation. *Proc of 24th Intl Conf on Machine Learning*, pages 265–272, 2007.

[15] K. Fukunaga and L. D. Hostetler. Optimization of k-nearest-neighbor density estimates. *IEEE Transactions on Information Theory*, 1973.

[16] E. Giné and D.M. Mason. Uniform in bandwidth estimation of integral functionals of the density function. *Scandinavian Journal of Statistics*, 35:739761, 2008.

[17] M. Goria, N. Leonenko, V. Mergel, and P. L. Novi Inverardi. A new class of random vector entropy estimators and its applications in testing statistical hypotheses. *Nonparametric Statistics*, 2004.

[18] Peter Hall and J. S. Marron. Estimation of integrated squared density derivatives. *Stat. Prob. Lett*, pages 109–115, 1987.

[19] A. O. Hero, J. Costa, and B. Ma. Asymptotic relations between minimal graphs and alpha-entropy. *Technical Report CSPL-334 Communications and Signal Processing Laboratory, The University of Michigan*, March 2003.

[20] A. O. Hero, B. Ma, O. Michel, and J. Gorman. Applications of entropic spanning graphs. *Signal Processing Magazine, IEEE*, 19(5):85 – 95, sep 2002.

[21] Marc M. Van Hulle. Edgeworth approximation of multivariate differential entropy. *Neural Computation*, 17(9):1903–1910, 2005.

[22] A. T. Ihler, J. W. Fisher III, and A. S. Willsky. Nonparametric hypothesis tests for statistical dependency. *IEEE Transactions on Signal Processing*, 52(8):2234–2249, August 2004.

[23] A.K. Jain. Image data compression: A review. *Proceedings of the IEEE*, 69(3):349 – 389, March 1981.

[24] A. Lakhina, M. Crovella, and C. Diot. Mining anomalies using traffic feature distributions. In *In ACM SIGCOMM*, pages 217–228, 2005.

[25] B. Laurent. Efficient estimation of integral functionals of a density. *The Annals of Statistics*, 24(2):659–681, 1996.

[26] N. Leonenko, L. Prozanto, and V. Savani. A class of rényi information estimators for multidimensional densities. *Annals of Statistics*, 36:2153–2182, 2008.

[27] N. Leonenko, L. Prozanto, and V. Savani. A class of rényi information estimators for multidimensional densities. *Annals of Statistics*, 36:2153–2182, 2008.

[28] E. Levina and P. J. Bickel. Maximum likelihood estimation of intrinsic dimension. In *Advances in Neural Information Processing Systems*, Cambridge, MA, 2005.

[29] E. Liitiäinen, A. Lendasse, and F. Corona. On the statistical estimation of rényi entropies. In *Proceedings of IEEE/MLSP 2009 International Workshop on Machine Learning for Signal Processing, Grenoble (France)*, September 2-4 2009.

[30] D. O. Loftsgaarden and C. P. Quesenberry. A nonparametric estimate of a multivariate density function. *Ann. Math. Statist.*, 1965.

[31] Y. P. Mack and M. Rosenblatt. Multivariate k-nearest neighbor density estimates. *Journal of Multivariate Analysis*, 9(1):1 – 15, 1979.

[32] E. G. Miller and J. W. Fisher III. ICA using spacings estimates of entropy. *Proc. 4th Intl. Symp. on ICA and BSS*, pages 1047–1052, 2003.

[33] D. S. Moore and J. W. Yackel. Consistency properties of nearest neighbor density function estimators. *The Annals of Statistics*, 1977.

[34] H. Neemuchwala and A. O. Hero. Image registration in high dimensional feature space. *Proc. of SPIE Conference on Electronic Imaging, San Jose*, January 2005.

[35] X. Nguyen, M. J. Wainwright, and M. I. Jordan. Estimating divergence functionals and the likelihood ratio by convex risk minimization. *Information Theory, IEEE Transactions on*, 56(11):5847 –5861, November 2010.

[36] D. Pál, B. Póczos, and C. Szepesvári. Estimation of R\'enyi Entropy and Mutual Information Based on Generalized Nearest-Neighbor Graphs. *ArXiv e-prints*, March 2010.

[37] B. Ranneby. The maximum spacing method. an estimation method related to the maximum likelihood method. *Scandinavian Journal of Statistics*, 11:93–112, 1984.

[38] V. C. Raykar and R. Duraiswami. Fast optimal bandwidth selection for kernel density estimation. In J. Ghosh, D. Lambert, D. Skillicorn, and J. Srivastava, editors, *Proceedings of the sixth SIAM International Conference on Data Mining*, pages 524–528, 2006.

[39] Xavier Saint Raymond. *Elementary Introduction to the Theory of Pseudodifferential Operators*. CRC Press, 1991.

[40] H. Singh, N. Misra, and V. Hnizdo. Nearest neighbor estimators of entropy. *The Annals of Statistics*, 2005.

[41] K. Sricharan, R. Raich, and A. O. Hero. Global performance prediction for divergence-based image registration criteria. In *Proc. IEEE Workshop on Statistical Signal Processing*, 2009.

[42] K. Sricharan, R. Raich, and A. O. Hero. Empirical estimation of entropy functionals with confidence. *ArXiv e-prints*, December 2010.

[43] B. van Es. Estimating functionals related to a density by class of statistics based on spacing. *Scandinavian Journal of Statistics*, 1992.

[44] O. Vasicek. A test for normality based on sample entropy. *Journal of the Royal Statistical Society. Series B (Methodological)*, 38:54–59, 1976.

[45] Q. Wang, S. R. Kulkarni, and S. Verdú. Divergence estimation of continuous distributions based on data-dependent partitions. *Information Theory, IEEE Transactions on*, 51(9):3064–3074, 2005.