

## Asymptotic Moments of Near Neighbor Distances for the Gaussian Distribution

Elia Liitiäinen\*

### Abstract

We study the moments  $E[d_{1,k}^\alpha]$  of the  $k$ -th nearest neighbor distance for independent identically distributed points in  $\mathfrak{R}^n$ . In the earlier literature, the case  $\alpha > n$  has been analyzed by assuming a bounded support for the underlying density. The boundedness assumption is removed by assuming the multivariate Gaussian distribution. In this case, the nearest neighbor distances show very different behavior in comparison to earlier results.

**Key words:** nearest neighbor; moments; gaussian; random geometry.

**AMS 2010 Subject Classification:** Primary 60D05.

Submitted to EJP on July 18, 2011, final version accepted November 11, 2011.

---

\*Department of Information and Computer Science, Aalto University School of Science and Technology, Espoo 021015, Finland

# 1 Introduction

Consider a set of independent identically distributed (i.i.d.) random variables  $(X_i)_{i=1}^M$  with a common density  $p(x)$  on  $\mathfrak{R}^n$ . We study the moments of the nearest neighbor distance

$$E[d_{1,k}^\alpha] \tag{1}$$

in the limit  $M \rightarrow \infty$ . The quantity (1) appears commonly in the literature on random geometric graphs, where directed and undirected nearest neighbor graphs are analyzed as special cases of more general frameworks [10, 11, 15]. In this paper, the nearest neighbor distance serves as the quantity of interest with the hope that in the future the ideas can be represented in a more abstract form.

The expectation (1) is also of interest in its own right and tends to appear under various scientific contexts. A significant application is found in the nonparametric estimation of Rényi entropies, where asymptotic analysis provides theoretically sound estimators [5, 7, 6, 9]. Moreover, nearest neighbor distances and distributions play a major role in the understanding of nonparametric estimation in general [1, 4, 13]. Finally, it should be mentioned that quantities related to (1) are encountered in physics, especially statistical mechanics and the theory of gases and liquids [3].

In the earlier literature, it has been shown that under general conditions ( $\Gamma$  denotes the Gamma function)

$$M^{\alpha/n} E[d_{1,k}^\alpha] \rightarrow V_n^{-\alpha/n} \frac{\Gamma(k + \alpha/n)}{\Gamma(k)} \int_{\mathfrak{R}^n} p(x)^{1-\alpha/n} dx$$

in the limit  $M \rightarrow \infty$  if  $0 < \alpha < n$  [14, 2]. However, the case  $\alpha > n$  is quite different and usually a boundedness condition must be imposed on the support of  $p(x)$ . As the contribution of this paper, we analyze what happens if  $\alpha > n$ , while  $p(x)$  is unbounded. To simplify matters, we examine only the multivariate Gaussian distribution  $p(x) = (2\pi)^{-n/2} e^{-\frac{1}{2}\|x\|^2}$  with the long term goal of extending the results to more general classes of densities. It turns out that the asymptotic behavior is very different to the case  $0 < \alpha < n$ . We show that if  $\alpha > n$ , then in the limit  $M \rightarrow \infty$ ,

$$(M \log^{\alpha/2+1-n} M) E[d_{1,k}^\alpha] \rightarrow \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!} \int_0^\infty g\left(\frac{1}{y}\right) dy,$$

where the definition of  $g$  depends on  $n, k$  and  $\alpha$  (see Section 3).

## 2 Definitions

We start with some basic definitions.  $V_n$  denotes the volume of the unit Euclidean ball in  $\mathfrak{R}^n$  and  $B(y, r)$  denotes the ball with center  $y$  and radius  $r$ .  $I(\cdot)$  refers to the indicator function of a random event. For a vector  $x \in \mathfrak{R}^n$ ,  $x^{(j)}$  denotes component  $j$  of that vector. The volume of a set  $A$  with respect to the Lebesgue measure is denoted by  $\lambda(A)$ . If  $g(r)$  is a function defined on an open subset of  $\mathfrak{R}$ , we denote the derivative of  $g$  by  $Dg$ .

$(X_i)_{i=1}^M$  is taken as an i.i.d. sample with  $X_i \in \mathfrak{R}^n$ . Each  $X_i$  follows a common density  $p(x)$ ; our work concerns the Gaussian case

$$p(x) = (2\pi)^{-n/2} e^{-\frac{1}{2}\|x\|^2}. \tag{2}$$

The first nearest neighbor of  $X_i$  is defined by (in the Euclidean norm, other norms are not considered in this paper)

$$N[i, 1] = \operatorname{argmin}_{1 \leq j \leq M, j \neq i} \|X_j - X_i\|$$

and by recursion, the  $k$ -th nearest neighbor is

$$N[i, k] = \operatorname{argmin}_{1 \leq j \leq M, j \notin \{i, N[i, 1], \dots, N[i, k-1]\}} \|X_j - X_i\|.$$

The corresponding  $k$ -th nearest neighbor distance is  $d_{i,k} = \|X_{N[i,k]} - X_i\|$ . The goal of the paper is to analyze

$$E[d_{i,k}^\alpha] \tag{3}$$

in the limit  $M \rightarrow \infty$  with everything else fixed. Because the sample is independent identically distributed (i.i.d), we set  $i = 1$ .

Throughout the paper there will be constants, which depend on some variables, but not on the others. Such variables are denoted by  $c(\dots)$ , where inside the parentheses we indicate the dependency. Strictly speaking,  $c$  is a function of some variables, but in the standard convention, it will be called a constant. During the course of our proofs, several different unknown constants will emerge. To keep them separate, lower indices (in the form  $c_i$ ) are used.

General error terms, which can be bounded but not written in closed form, will be denoted by  $R$  (or  $R_i$  with a lower index  $i$ ). After the appearance of each such term, we write an equation of the form

$$|R| \leq c(\dots)f(\dots),$$

where  $c$  is a constant and  $f$  is a function of  $M$  or some other variables. Inside proofs, the Big-Oh notation will be invoked as another way to express unknown but negligible terms.

### 3 Main Results and Previous Work

The analysis of nearest neighbor distances can be viewed as part of the general framework of random geometric graphs. In this field, results are established for quantities of the form  $\xi(X_1, (X_i)_{i=1}^M)$ , where  $\xi$  has some locality properties. By imposing higher levels of abstraction, very general functions can be analyzed as long as locality arguments are available. We refer to [10, 11, 15] as a starting point to understand the issues arising in the field.

However, abstract theories do not directly give exact information about the asymptotic behavior of the moments (3). The step towards concretizing the results concerning nearest neighbor graphs was taken in [14, 12]. The following has been proven:

**Theorem 1.** *Suppose that  $0 < \alpha < n$  and  $p(x)$  is a density with*

$$\int_{\mathfrak{R}^n} p(x)^{1-\alpha/n} dx < \infty$$

and

$$\int_{\mathfrak{R}^n} \|x\|^r p(x) dx < \infty$$

for some  $r > \alpha n / (n - \alpha)$ . Then

$$M^{\alpha/n} E[d_{1,k}^\alpha] \rightarrow V_n^{-\alpha/n} \frac{\Gamma(k + \alpha/n)}{\Gamma(k)} \int_{\mathbb{R}^n} p(x)^{1-\alpha/n} dx$$

in the limit  $M \rightarrow \infty$ .  $\Gamma(\cdot)$  refers to the Gamma function. If  $\alpha \geq n$ , the limit holds if  $p(x)$  is bounded from below and above on a bounded convex set  $\mathcal{X}$  with  $p(x) = 0$  when  $x \notin \mathcal{X}$ .

As a downside, Theorem 1 imposes the convexity requirement on  $\mathcal{X}$  if  $\alpha > n$ . Furthermore, it does not provide a rate of convergence. These issues have been addressed by the concrete approach in [2], where it was shown that if  $\inf_{x \in \mathcal{X}} p(x) > 0$  and  $p(x)$  has a bounded gradient on  $\mathcal{X}$ , then under rather weak conditions on the space  $\mathcal{X}$ , we have

$$M^{\alpha/n} E[d_{1,k}^\alpha] = V_n^{-\alpha/n} \frac{\Gamma(k + \alpha/n)}{\Gamma(k)} \int_{\mathcal{X}} p(x)^{1-\alpha/n} dx + O(M^{-1/n+\rho})$$

for any  $\rho > 0$  removing the convexity requirement.

As a common factor between the results, observe that in the case  $\alpha > n$ , two requirements must be satisfied:

1. The set  $\mathcal{X}$  must be bounded.
2.  $\inf_{x \in \mathcal{X}} p(x) > 0$ .

In this paper we ask, what happens when neither 1. nor 2. hold but  $\alpha > n$  (the case  $\alpha = n$  is not addressed). The early works in random geometry took the uniform distributions as a case of special interest. Analogously, we choose the Gaussian density (2) as our target of study.

It turns out that the behavior for  $\alpha > n$  is very different to Theorem 1 for the Gaussian distribution. As the main contribution of the paper, we prove the following.

**Theorem 2.** Suppose that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then

$$(M \log^{\alpha/2+1-n} M) E[d_{1,k}^\alpha] \rightarrow \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!} \int_0^\infty g\left(\frac{1}{y}\right) dy \tag{4}$$

in the limit  $M \rightarrow \infty$  with

$$g(t) = \int_0^\infty \omega^{k-1} e^{-\omega} f^{-1}(\omega t)^\alpha d\omega,$$

where  $f^{-1}$  refers to the inverse function of

$$f(t) = t^n \int_{B(0,1)} e^{ty^{(1)}} dy.$$

Moreover, the function  $g(y^{-1})$  is integrable on  $(0, \infty)$  and consequently the limit (4) is finite.

The main difference to Theorem 1 is that now  $E[d_{1,k}^\alpha]$  is of order  $M^{-1}(\log M)^{n-\alpha/2-1}$  instead of  $M^{-\alpha/n}$ . Theorem 2 can be further developed by analyzing the rate of convergence. In fact, the results suggest that even in the well studied case  $0 < \alpha < n$  the rates of convergence obtained for example in [2, 8] do not hold in the unbounded case especially when  $\alpha$  is close to  $n$ . The rather deep questions related to the rates of convergence are left as an important topic of future research.

Another open question is the extension to a general density  $p$ , which the author believes is possible. This could possibly unify the case with boundary effect [8] and the more general unbounded case.

## 4 Outline of the Proof

We will use the small ball probability

$$\omega_x(r) = \int_{B(x,r)} p(y) dy$$

due to its useful distribution free properties. In fact, [13, 2] shows that the distribution of the quantity  $\omega_{X_1}(d_{1,k})$  does not depend on the density  $p$  and moreover, concentrates on values of order  $M^{-1}$ . Another useful fact is that conditioning on  $X_1$  does not change the distribution of  $\omega_{X_1}(d_{1,k})$ . We approximate

$$\begin{aligned} \omega_x(r) &= (2\pi)^{-n/2} \int_{B(x,r)} e^{-\frac{1}{2}\|y\|^2} dy \\ &= (2\pi)^{-n/2} \int_{B(x,r)} e^{-\frac{1}{2}\|x\|^2 - x^T(y-x) - \frac{1}{2}\|y-x\|^2} dy \\ &\approx p(x) \int_{B(0,r)} e^{-x^T y} dy = p(x)r^n \int_{B(0,1)} e^{-rx^T y} dy \end{aligned} \quad (5)$$

assuming that  $e^{-\frac{1}{2}r^2}$  is close to 1. By a change of variables (rotation inside the last integral in (5)) we have

$$\omega_x(r) \approx p(x)r^n \int_{B(0,1)} e^{-r\|x\|y^{(1)}} dy.$$

Now if we take  $f(t) = t^n \int_{B(0,1)} e^{-ty^{(1)}} dy$ , then  $\|x\|^n \omega_x(r) \approx p(x)f(\|x\|r)$  and we solve

$$r \approx \frac{f^{-1}\left(\frac{\|x\|^n \omega_x(r)}{p(x)}\right)}{\|x\|}.$$

$f^{-1}$  refers to the inverse of  $f$ . By substituting  $d_{1,k}$  in place of  $r$  and  $\|X_1\|$  in place of  $\|x\|$ , we get conditionally on  $X_1$

$$E[d_{1,k}^\alpha] \approx E \left[ E \left[ \frac{f^{-1}\left(\frac{\|X_1\|^n \omega_{X_1}(d_{1,k})}{p(X_1)}\right)^\alpha}{\|X_1\|^\alpha} \middle| X_1 \right] \right].$$

The argument for  $f^{-1}$  looks rather complicated. However, because the conditional distribution of  $\omega_{X_1}(d_{1,k})$  does not depend on the density  $p(x)$  or  $X_1$ , it would be sufficient to somehow control the dependency on  $X_1$ . Our strategy can be summarized as dividing  $\mathfrak{R}^n$  into the three regions  $S_1$ ,  $S_2$  and  $S_3$  together with decomposing

$$\begin{aligned} E[d_{1,k}^\alpha] &= \int_{S_1} E[d_{1,k}^\alpha | X_1 = x] p(x) dx + \int_{S_2} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \\ &\quad + \int_{S_3} E[d_{1,k}^\alpha | X_1 = x] p(x) dx. \end{aligned}$$

The three sets depend on a variable  $0 < \epsilon < 1$  and the number of samples  $M$ . We think  $\epsilon > 0$  as a parameter, which at the end of the analysis is set to approach zero after first taking the limit  $M \rightarrow \infty$ . As a sidenote, it should be clear at this point that the parameters  $(n, k, \alpha)$  are assumed to stay fixed all the time.

The motivation for  $S_1$  might be seen in the idea of performing a Taylor expansion of  $f^{-1}(\cdot)^\alpha$  at zero, which might render the analysis into the well-known case [2]. Keeping in mind that  $\omega_{X_1}(d_{1,k})$  is of order of magnitude  $M^{-1}$ , we take (the definition applies for any  $n \geq 1$ )

$$\begin{aligned} S_1 &= \{x \in \mathfrak{R}^n : p(x) > \frac{\log^{n/2} M}{\epsilon M}\} \\ &= \{x \in \mathfrak{R}^n : \|x\| < \sqrt{2 \log M - n \log \log M + 2 \log \epsilon - n \log 2\pi}\}; \end{aligned} \quad (6)$$

then for large  $M$ ,  $\|X_1\| = O(\sqrt{\log M})$  when  $X_1 \in S_1$  and

$$\frac{\|X_1\|^n \omega_{X_1}(d_{1,k})}{p(X_1)} = O(\epsilon)$$

by substituting  $\omega_{X_1}(d_{1,k}) = \frac{1}{M}$  to analyze the order of magnitude. If  $\epsilon$  is small, then this shows that the argument of  $f^{-1}$  is small suggesting that a Taylor expansion might be possible. However, during the course of the proof, it turns out that points in  $S_1$  contribute little in comparison to the set

$$S_2 = \{x \in \mathfrak{R}^n : \frac{\epsilon \log^{n/2} M}{M} \leq p(x) \leq \frac{\log^{n/2} M}{\epsilon M}\}. \quad (7)$$

In this case, a Taylor expansion does not seem possible. Fortunately, we are able to show that conditionally on  $X_1 \in S_2$ , the variable

$$Y = \frac{Mp(X_1)}{\log^{n/2} M} \quad (8)$$

is approximately uniformly distributed on  $[\epsilon, \epsilon^{-1}]$  and moreover, it is independent of  $\omega_{X_1}(d_{1,k})$ . This is useful, because for large  $M$ ,  $\|X_1\| \approx \sqrt{2 \log M}$  and we get

$$E[d_{1,k}^\alpha | X_1 \in S_2] \approx E \left[ \frac{f^{-1} \left( \frac{2^{n/2} \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha}{(2 \log M)^{\alpha/2}} | X_1 \in S_2 \right]. \quad (9)$$

Because the probability  $P(X_1 \in S_2)$  turns out to admit a convenient asymptotic expression, it is possible to use Equation (9) to estimate the quantity

$$\int_{S_2} E[d_{1,k}^\alpha | X_1 = x] p(x) dx = E[d_{1,k}^\alpha | X_1 \in S_2] P(X_1 \in S_2).$$

In addition to  $S_1$  and  $S_2$ , there is the set

$$S_3 = \{x \in \mathfrak{R}^n : p(x) < \frac{\epsilon \log^{n/2} M}{M}\}. \quad (10)$$

However, similarly as  $S_1$ , nearest neighbor distances corresponding to  $X_1 \in S_3$  turn out to have a negligible effect if  $\epsilon$  is small.

## 5 Auxiliary Results

In this section, we give some results and applications for  $\omega_{X_1}(d_{1,k})$ , where

$$\omega_x(r) = \int_{B(x,r)} p(x) dx.$$

The following result characterizes the distribution of  $\omega_{X_1}(d_{1,k})$ , which conveniently does not depend on  $X_1$  or the density  $p(x)$ .

**Lemma 1.** *Given  $X_1$ , the conditional density of  $\omega_{X_1}(d_{1,k})$  is given by*

$$p_\omega(\omega|X_1) = p_\omega(\omega) = k \binom{M-1}{k} \omega^{k-1} (1-\omega)^{M-k-1}. \quad (11)$$

Moreover,

$$E[\omega_{X_1}(d_{1,k})^\alpha | X_1] = \frac{\Gamma(k + \alpha/n) \Gamma(M)}{\Gamma(k) \Gamma(M + \alpha/n)}. \quad (12)$$

*Proof.* Equation (11) can be derived from the the cumulative distribution function in Equation (4.35) of [2]. Some algebraic manipulation is needed to simplify the first derivative of the sum of terms appearing in [2] in order to reach the simpler formula (11).  $\square$

It is useful to observe that for any  $\beta > 0$ ,

$$\frac{\Gamma(M + \beta)}{\Gamma(M)} = M^\beta + O(M^{\beta-1}) \quad (13)$$

to understand better the moments (12). The following lemma is useful for technical reasons.

**Lemma 2.** *Assume that  $p(x)$  is the multivariate Gaussian distribution (2). Then for  $0 < r < 1$  and  $x \in \mathfrak{R}^n$ ,*

$$\omega_x(r) \geq c p(x) r^n$$

for some constant  $c(n) > 0$ .

*Proof.* By a slight modification to Equation (5), we have

$$\omega_x(r) \geq p(x) e^{-\frac{1}{2}r^2} \int_{B(0,r)} e^{-x^T y} dy \geq \frac{e^{-\frac{1}{2}}}{2} V_n p(x) r^n. \quad (14)$$

$\square$

The moments  $E[d_{1,k}^\alpha | X_1]$  do not get too large if  $\|X_1\|$  stays close enough to the origin. The following lemma can be proven for example by observing that  $d_{1,k}^\alpha \leq \sum_{i=2}^{k+1} \|X_1 - X_i\|^\alpha$ :

**Lemma 3.** Assume that  $p(x)$  is the multivariate Gaussian distribution (2). Then for  $x \in \mathfrak{R}^n$ ,  $M > 2k$  and  $\alpha > 0$

$$E[d_{1,k}^\alpha | X_1 = x] \leq c(\|x\|^\alpha + 1)$$

for some constant  $c(n, k, \alpha)$ .

Next we show that the  $\alpha$ -moments are at most of order  $(p(x)M)^{-\alpha/n}$  if the quantity inside the parentheses does not get too small. The result is an application of Lemmas 1-2. Without losing generality, we prove the claim after some threshold  $M_0$ , which is natural as in any case later the limit  $M \rightarrow \infty$  is taken. As a somewhat subtle detail, we will generally adopt this way of expressing our statements in those cases, where proving the claim for all  $M > 0$  is not an obvious task.

**Lemma 4.** Suppose that  $p(x)$  is the multivariate Gaussian distribution (2) and fix any  $\delta > 0$ . Then if  $p(x) > \frac{\delta \log^{n/2} M}{M}$ , we find a threshold  $M_0(n, k, \alpha, \delta)$  such that for all  $M > M_0$ , we have

$$E[d_{1,k}^\alpha | X_1 = x] \leq c(p(x)M)^{-\alpha/n}$$

for some constant  $c(n, k, \alpha)$ .

*Proof.* We decompose

$$E[d_{1,k}^\alpha | X_1 = x] = E[d_{1,k}^\alpha I(d_{1,k} \leq 1) | X_1 = x] + E[d_{1,k}^\alpha I(d_{1,k} > 1) | X_1 = x]. \quad (15)$$

We consider next the first term in the right side. By Lemma 2,

$$\frac{d_{1,k}^n}{\omega_{X_1}(d_{1,k})} I(d_{1,k} \leq 1) \leq \frac{c_1}{p(X_1)} \quad (16)$$

(for some constant  $c_1(n)$ ) and using this we have by Lemma 1 together with Equations (13) and (16),

$$E[d_{1,k}^\alpha I(d_{1,k} \leq 1) | X_1 = x] \leq \frac{c_1^{\alpha/n} E[\omega_{X_1}(d_{1,k})^{\alpha/n} | X_1 = x]}{p(x)^{\alpha/n}} \leq \frac{c_2}{(p(x)M)^{\alpha/n}} \quad (17)$$

for some constant  $c_2(n, k, \alpha)$ . We have proven the claim for the first term in (15). For the second term, we apply Hölder's inequality:

$$E[d_{1,k}^\alpha I(d_{1,k} > 1) | X_1 = x] \leq \sqrt{P(d_{1,k} > 1 | X_1 = x)} \sqrt{E[d_{1,k}^{2\alpha} | X_1 = x]}. \quad (18)$$

$\omega_x(r)$  is a strictly increasing function with respect to  $r$  and Equation (16) implies that  $\omega_x(1) \geq c_1^{-1} p(x)$ . Using this fact, integration by parts and the inequalities  $k \binom{M-1}{k} \leq M^k$  and  $1 - \omega \leq e^{-\omega}$  together with Lemma 1 and a change of variables, we have

$$\begin{aligned} P(d_{1,k} > 1 | X_1 = x) &= P(\omega_{X_1}(d_{1,k}) > \omega_{X_1}(1) | X_1 = x) \\ &\leq \left( \frac{M}{M-k-1} \right)^k \int_{c_1^{-1}(M-k-1)p(x)}^{M-k-1} \omega^{k-1} e^{-\omega} d\omega \\ &\leq k! \left( \frac{M}{M-k-1} \right)^k (c_1^{-1} M p(x) + 1)^k e^{-c_1^{-1}(M-k-1)p(x)}. \end{aligned} \quad (19)$$



The second line in (19) can be easily calculated in closed form, but for our purposes it is convenient to use the upper bound to simplify the notation. Assuming  $M > 2k + 2$ , we have

$$P(d_{1,k} > 1 | X_1 = x) \leq 2^k k! (c_1^{-1} p(x) M + 1)^k e^{-\frac{1}{2} c_1^{-1} p(x) M} \leq c_3 e^{-\frac{1}{4} c_1^{-1} p(x) M} \quad (20)$$

for some  $c_3(n, k)$ . By the assumption  $p(x) > \frac{\delta \log^{n/2} M}{M}$

$$\|x\| \leq \sqrt{2 \log M - n \log \log M - 2 \log \delta - n \log(2\pi)} \leq \sqrt{3 \log M}$$

after some threshold  $M_0(n, \delta)$  and for all  $M > M_0$ . By Lemma 3 we then have

$$E[d_{1,k}^{2\alpha} | X_1 = x] \leq c_4 \log^\alpha M \quad (21)$$

for some constant  $c_4(n, k, \alpha)$  (assuming trivially  $M > 1$ ). Equations (20) and (21) together with (18) now imply

$$E[d_{1,k}^\alpha I(d_{1,k} > 1) | X_1 = x] \leq \sqrt{c_3 c_4} e^{-\frac{1}{8} c_1^{-1} p(x) M} \log^{\alpha/2} M. \quad (22)$$

The assumption  $p(x) M \geq \delta \log^{n/2} M$  implies that Equation (22) approaches zero faster than  $(p(x) M)^{-\alpha/n}$  in the limit  $M \rightarrow \infty$ .  $\square$

We formalize the argument in Section 4, which connects  $\omega_x(r)$  to the function  $f$ :

**Lemma 5.** *Suppose that  $p(x)$  is the multivariate Gaussian distribution (2). Then*

$$\|x\|^n \omega_x(r) = p(x) f(\|x\| r) - R$$

with

$$f(t) = t^n \int_{B(0,1)} e^{ty^{(1)}} dy$$

and  $0 \leq R \leq p(x) r^2 f(\|x\| r)$ .  $f$  is defined and continuous on  $[0, \infty)$  and it has the range  $[0, \infty)$ . It is also strictly increasing implying the existence of an inverse function  $f^{-1} : [0, \infty) \mapsto [0, \infty)$ .

*Proof.* The proof involves extracting the error term and bounding it. By rearranging terms and a change of variables (see also Equation (5))

$$\begin{aligned} \|x\|^n \omega_x(r) &= (2\pi)^{-n/2} \|x\|^n \int_{B(x,r)} e^{-\frac{1}{2} \|y\|^2} dy \\ &= p(x) (\|x\| r)^n \int_{B(0,1)} e^{rx^T y} dy - A \end{aligned} \quad (23)$$

with

$$A = p(x) (\|x\| r)^n \int_{B(0,1)} e^{rx^T y} (1 - e^{-\frac{1}{2} r^2 \|y\|^2}) dy. \quad (24)$$

The main task is to bound  $A$ . This is achieved by the mean-value theorem: for  $\|y\| \leq 1$  and  $r > 0$ ,

$$1 - e^{-\frac{1}{2} r^2 \|y\|^2} = \frac{1}{2} r^2 \|y\|^2 e^{-\delta} \leq r^2$$

for some  $\delta \in (0, \infty)$ . This inequality implies that

$$\begin{aligned} 0 \leq A &\leq p(x)(\|x\|r)^n r^2 \int_{B(0,1)} e^{rx^T y} dy \\ &\leq p(x)(\|x\|r)^n r^2 \int_{B(0,1)} e^{r\|x\|y^{(1)}} dy = p(x)r^2 f(\|x\|r). \end{aligned}$$

In the last inequality, the vectors have been conveniently rotated. The same rotation shows that in (23), we have

$$p(x)(\|x\|r)^n \int_{B(0,1)} e^{rx^T y} dy = p(x)f(\|x\|r).$$

□

For  $t > 0$ , we define

$$g(t) = \int_0^\infty \omega^{k-1} e^{-\omega} f^{-1}(\omega t)^\alpha d\omega. \quad (25)$$

The integral always exists because  $f^{-1}$  is a non-negative function. We show that  $g$  approaches zero at least as fast as  $t^{\alpha/n}$  and grows at most logarithmically if  $t \rightarrow \infty$ . The same holds for  $f^{-1}(t)^\alpha$ :

**Lemma 6.** *The functions  $f(t)$  and  $g(t)$  are bounded by*

$$0 \leq g(t) + f^{-1}(t)^\alpha \leq ct^{\alpha/n}$$

on  $(0, 1]$  for some constant  $c(n, k, \alpha)$ . On  $(1, \infty)$  we have

$$0 \leq g(t) + f^{-1}(t)^\alpha \leq c(1 + \log^\alpha t).$$

*Proof.* 1. Bounds on  $f^{-1}$

Consider  $t \in (0, 1)$ . For any  $z > 2V_n^{-1/n}t^{1/n}$ , we have

$$f(z) > \frac{2t \int_{B(0,1)} e^{zy^{(1)}} dy}{V_n} > t.$$

This implies that  $f^{-1}(t) \leq 2V_n^{-1/n}t^{1/n}$ . Next assume that  $t > 1$ . Take  $z > 2\log t + A + 1$  with

$$A = \lambda(B(0, 1) \cap \{x \in \mathfrak{R}^n : x^{(1)} > \frac{1}{2}\})^{-1}.$$

Then

$$f(z) > A \int_{B(0,1) \cap \{x: x^{(1)} > \frac{1}{2}\}} e^{2y^{(1)} \log t} dy > t.$$

This means that  $f^{-1}(t) \leq 2 \log t + A + 1$ . The outcome for  $f^{-1}(t)^\alpha$  follows by recalling that  $(a+b)^\alpha \leq 2^\alpha(a^\alpha + b^\alpha)$  for any  $a, b > 0$ .

## 2. The function $g$

Bounds for  $g$  can be established for example by using

$$g(t) \leq \sum_{i=1}^{\infty} 2^{i(k-1)} e^{-2^{i-1}} f^{-1}(2^i t).$$

When  $t \in (0, 1)$  the proof can be established by examining the terms with  $2^i t < 1$  and  $2^i t \geq 1$  separately, whereas for  $t > 1$  a straightforward application of the logarithmic upper bound for  $f^{-1}$  gives the result.  $\square$

## 6 Region $S_1$

Recall that region  $S_1$  is defined by

$$\begin{aligned} S_1 &= \{x \in \mathfrak{R}^n : p(x) > \frac{\log^{n/2} M}{\epsilon M}\} \\ &= \{x \in \mathfrak{R}^n : \|x\| < \sqrt{2 \log M - n \log \log M + 2 \log \epsilon - n \log(2\pi)}\}. \end{aligned} \quad (26)$$

It may happen that  $S_1$  is an empty set; from now on we always assume that  $M$  is large enough in comparison to  $\epsilon^{-1}$  and  $n$  in order to ensure that  $S_1$  is non-empty with a positive volume. Similar convention is adopted for the sets  $S_2$  and  $S_3$ .

As stated in Section 4,  $0 < \epsilon < 1$  is a fixed constant until the end, where the limit  $\epsilon \rightarrow 0$  is taken after the limit  $M \rightarrow \infty$ . We define (assuming that  $\alpha > n$ )

$$i^* = \lceil \log^{-1} 2 \frac{n \log \log M}{\alpha - n} \rceil + 1.$$

$\lceil \cdot \rceil$  refers to the integer part of the number inside the bracket. As our proof strategy,  $S_1$  is divided into smaller subsets, which are easier to control with the tools we have available this far:

$$\begin{aligned} \tilde{S}_{1,i} &= \{x \in \mathfrak{R}^n : 2^i \frac{\log^{n/2} M}{\epsilon M} < p(x) \leq 2^{i+1} \frac{\log^{n/2} M}{\epsilon M}\} \\ &= \{x \in \mathfrak{R}^n : \|x\| \in [a_i, b_i)\} \end{aligned} \quad (27)$$

( $0 \leq i \leq i^*$ ) with

$$\begin{aligned} a_i &= \sqrt{2 \log M - n \log \log M - 2(i+1) \log 2 + 2 \log \epsilon - n \log(2\pi)} \\ b_i &= \sqrt{2 \log M - n \log \log M - 2i \log 2 + 2 \log \epsilon - n \log(2\pi)}. \end{aligned}$$

The remaining part is denoted by

$$S_{1,C} = S_1 \setminus \cup_{i=0}^{i^*} \tilde{S}_{1,i}.$$

The following bounds the nearest neighbor distance when  $X_1 \in \tilde{S}_{1,i}$ .

**Lemma 7.** Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then there exists a threshold  $M_0(n, k, \alpha, \epsilon) > 0$  such that for all  $M > M_0$  and  $0 \leq i \leq i^*$ ,

$$\int_{\tilde{S}_{1,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq 2^{i(1-\alpha/n)} c \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M}$$

for some constant  $c(n, k, \alpha)$ .

*Proof.* By Lemma 4,

$$\begin{aligned} \int_{\tilde{S}_{1,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx &\leq c_1 M^{-\alpha/n} \int_{\tilde{S}_{1,i}} p(x)^{1-\alpha/n} dx \\ &\leq 2^{i(1-\alpha/n)} c_1 \epsilon^{\alpha/n-1} \frac{\log^{n/2-\alpha/2} M}{M} \lambda(\tilde{S}_{1,i}) \end{aligned} \quad (28)$$

for some constant  $c_1(n, k, \alpha)$  and  $M_0(n, k, \alpha, \epsilon)$ . We should now compute the volume  $\lambda(\tilde{S}_{1,i})$ . The set  $\tilde{S}_{1,i}$  consists of points  $x \in \mathfrak{R}^n$  with  $\|x\|$  in the interval  $[a_i, b_i)$  and the volume of the set is  $\lambda(\tilde{S}_{1,i}) = V_n(b_i^n - a_i^n)$ . By a Taylor expansion,

$$\begin{aligned} a_i^n &= 2^{n/2} \log^{n/2} M \left( 1 - \frac{n^2 \log \log M + 2n(i+1) \log 2 - 2n \log \epsilon + n^2 \log(2\pi)}{4 \log M} \right) \\ &\quad + R_1 \\ b_i^n &= 2^{n/2} \log^{n/2} M \left( 1 - \frac{n^2 \log \log M + 2ni \log 2 - 2n \log \epsilon + n^2 \log(2\pi)}{4 \log M} \right) + R_2 \end{aligned} \quad (29)$$

in the limit  $M \rightarrow \infty$  with everything else fixed and

$$|R_1| + |R_2| \leq c_2 \frac{\log^2 \log M}{\log^{2-n/2} M}$$

with  $c_2(n, k, \alpha, \epsilon)$  independent of  $i$ . The Taylor expansions imply that  $\lambda(\tilde{S}_{1,i}) = V_n(b_i^n - a_i^n) \leq c_3 \log^{n/2-1} M$  for some constant  $c_3(n, k, \alpha, \epsilon)$ . By substitution into (28), we have

$$\int_{\tilde{S}_{1,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq 2^{i(1-\alpha/n)} c_1 c_3 \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M},$$

where the bound holds for  $0 \leq i \leq i^*$ . □

After removing the sets  $\tilde{S}_{1,i}$ , we are left with  $\tilde{S}_{1,C}$ . However, it does not pose problems.

**Lemma 8.** Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then there exists a threshold  $M_0(n, k, \alpha, \epsilon)$  such that for any  $M > M_0$ , we have

$$\int_{\tilde{S}_{1,C}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M}$$

for some constant  $c(n, k, \alpha)$ .

*Proof.* By Lemma 4 and the definition of  $\tilde{S}_{1,C}$ ,

$$\begin{aligned} \int_{\tilde{S}_{1,C}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx &\leq c_1 M^{-\alpha/n} \int_{\tilde{S}_{1,C}} p(x)^{1-\alpha/n} dx \\ &\leq 2^{i^*(1-\alpha/n)} c_1 \epsilon^{\alpha/n-1} \frac{\log^{n/2-\alpha/2} M}{M} \lambda(\tilde{S}_{1,C}) \end{aligned} \quad (30)$$

for some constant  $c_1(n, k, \alpha)$ . It is a simple task to show that for all  $x \in \tilde{S}_{1,C}$  we have  $\|x\| \leq \sqrt{3 \log M}$  once  $M$  exceeds some threshold  $M_0(n, k, \alpha, \epsilon)$ . This implies that

$$\lambda(\tilde{S}_{1,C}) \leq 3^{n/2} V_n \log^{n/2} M. \quad (31)$$

Substituting Equation (31) and the inequality  $2^{i^*(1-\alpha/n)} \leq \log^{-1} M$  into (30) yields

$$\int_{\tilde{S}_{1,C}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq 3^{n/2} c_1 V_n \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M}.$$

□

Lemmas 7 and 8 imply that for  $\alpha > n$  and  $M > M_0$

$$\int_{S_1} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M} \left(1 + \sum_{i=0}^{i^*} 2^{i(1-\alpha/n)}\right) \quad (32)$$

for some constant  $c(n, k, \alpha)$ . We conclude

**Lemma 9.** *Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then there exists a threshold  $M_0(n, k, \alpha, \epsilon)$  such that for any  $M > M_0$ , we have*

$$\int_{S_1} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c \epsilon^{\alpha/n-1} \frac{\log^{n-\alpha/2-1} M}{M}$$

for some constant  $c(n, k, \alpha)$ .

## 7 Region $S_2$

Region 2 is defined by

$$S_2 = \{x \in \mathfrak{R}^n : \frac{\epsilon \log^{n/2} M}{M} \leq p(x) \leq \frac{\log^{n/2} M}{\epsilon M}\}. \quad (33)$$

Again,  $M$  is assumed to be large enough to ensure that  $S_2$  has a positive volume. It is necessary to obtain an approximation to  $P(X_1 \in S_2)$ . This can be done rather straightforwardly:

**Lemma 10.** Assuming that  $p(x)$  is the multivariate Gaussian distribution (2), it holds that

$$P(X_1 \in S_2) = \frac{2^{n/2-1} n V_n \log^{n-1} M}{\epsilon M} (1 - \epsilon^2) + R,$$

where for some constant  $c(n, \epsilon)$ ,

$$|R| \leq c \frac{\log^2 \log M \log^{n-2} M}{M}.$$

*Proof.*  $S_2$  consists of points  $x \in \mathfrak{R}^n$  with  $\|x\| \in [a, b]$  and

$$a = \sqrt{2 \log M - n \log \log M + 2 \log \epsilon - n \log(2\pi)} \quad (34)$$

$$b = \sqrt{2 \log M - n \log \log M - 2 \log \epsilon - n \log(2\pi)}. \quad (35)$$

By some algebraic manipulation,

$$P(X_1 \in S_2) = (2\pi)^{-n/2} n V_n \int_a^b x^{n-1} e^{-\frac{1}{2}x^2} dx = I_1 + I_2 + I_3$$

with

$$I_1 = \frac{n V_n \log^{n/2} M}{\epsilon M} a^{n-1} \int_a^b e^{-(x-a)a} dx = \frac{n V_n \log^{n/2} M}{\epsilon M} a^{n-2} (1 - e^{-(b-a)a})$$

$$I_2 = \frac{n V_n \log^{n/2} M}{\epsilon M} \int_a^b (x^{n-1} - a^{n-1}) e^{-(x-a)a - \frac{1}{2}(x-a)^2} dx$$

$$I_3 = \frac{n V_n \log^{n/2} M}{\epsilon M} a^{n-1} \int_a^b (e^{-\frac{1}{2}(x-a)^2} - 1) e^{-(x-a)a} dx.$$

During the proof it is easiest to employ the Big-Oh notation. Such error terms depend here on  $n$  and  $\epsilon$ .

### 1. The term $I_1$

By a Taylor expansion (see also Equation (29)), it can be shown that

$$b - a = \frac{\sqrt{2} \log \epsilon^{-1}}{\sqrt{\log M}} + O\left(\frac{\log^2 \log M}{\log^{3/2} M}\right) \quad (36)$$

and for any  $\beta > 0$ ,

$$a^\beta = 2^{\beta/2} \log^{\beta/2} M + O\left(\frac{\log \log M}{\log^{1-\beta/2} M}\right). \quad (37)$$

Using (36) and (37) with  $\beta = 1$ , we have

$$1 - e^{-(b-a)a} = 1 - e^{2 \log \epsilon + O\left(\frac{\log^2 \log M}{\log M}\right)} = 1 - \epsilon^2 + O\left(\frac{\log^2 \log M}{\log M}\right). \quad (38)$$

The remainder terms depend on  $n, \epsilon$  and  $M$ . Using Equations (38) and (37) with  $\beta = n - 2$  in the expression for  $I_1$  yields

$$I_1 = \frac{2^{n/2-1} n V_n \log^{n-1} M}{\epsilon M} (1 - \epsilon^2) + O\left(\frac{\log^2 \log M \log^{n-2} M}{M}\right).$$

### 2. The term $I_2$

By the mean value theorem, for some constant  $c_1(n, \epsilon)$  we have  $|x^{n-1} - a^{n-1}| \leq c_1 \log^{n/2-3/2} M$  for  $x \in [a, b]$ . Also,  $a^{-1} \leq c_2 \log^{-1/2} M$  for some  $c_2(n, \epsilon)$ . We have

$$I_2 \leq \frac{c_1 n V_n \log^{n-3/2} M}{\epsilon M} \int_a^b e^{-(x-a)a} dx \leq \frac{c_1 c_2 n V_n \log^{n-2} M}{\epsilon M}.$$

### 3. The term $I_3$

Now by  $e^{-\frac{1}{2}(x-a)^2} \geq e^{-\frac{1}{2}(b-a)^2}$  (for  $x \in [a, b]$ ), we have

$$|I_3| \leq \frac{n V_n \log^{n/2} M}{\epsilon M} a^{n-2} (1 - e^{-\frac{1}{2}(b-a)^2}). \quad (39)$$

Moreover, by the expansion for  $b - a$  appearing in Equation (36),

$$1 - e^{-\frac{1}{2}(b-a)^2} = \frac{1}{2}(b-a)^2 + O((b-a)^4) \leq \frac{c_3}{\log M} \quad (40)$$

for some constant  $c_3(n, \epsilon)$ . Finally,

$$a^{n-2} \leq c_4 \log^{n/2-1} M. \quad (41)$$

for some constant  $c_4(n, \epsilon)$ . Substituting (40)-(41) into (39) yields

$$|I_3| \leq \frac{c_3 c_4 n V_n \log^{n-2} M}{\epsilon M}.$$

The proof is finished since the terms  $I_1, I_2$  and  $I_3$  have been addressed. □

In general, to establish asymptotics, it is useful to truncate  $d_{1,k}$  to avoid too large values. To this end, we choose some  $L > 0$  (recall that at this point,  $\alpha, n, k$  and  $\epsilon$  stay fixed) and define the indicator

$$I_L = I\left(d_{1,k} < \frac{L}{\epsilon^{1/n} \sqrt{\log M}}\right).$$

The power for  $\log M$  is chosen to ensure the correct order of magnitude with large  $L$  rendering the event  $1 - I_L$  negligible. The following lemma verifies this fact; the bound is designed to hold after some threshold  $M_0$ , which depends on  $L$  itself. However, after the threshold we get an upper bound which goes exponentially to zero with respect to  $L$ .

**Lemma 11.** Suppose that  $p(x)$  is the multivariate Gaussian distribution (2). Then for any  $L > 0$ , there exists a threshold  $M_0(n, k, \alpha, \epsilon, L)$  and a positive constant  $c(n, k, \alpha, \epsilon)$  such that for all  $M > M_0$ , it holds that

$$E[d_{1,k}^\alpha (1 - I_L) | X_1 \in S_2] \leq ce^{-c^{-1}L^n} \log^{-\alpha/2} M.$$

*Proof.* The proof employs Hölder's inequality:

$$E[d_{1,k}^\alpha (1 - I_L) | X_1 \in S_2] \leq \sqrt{E[d_{1,k}^{2\alpha} | X_1 \in S_2] P(d_{1,k} > \frac{L}{\epsilon^{1/n} \sqrt{\log M}} | X_1 \in S_2)}. \quad (42)$$

By Lemma 4 and the definition of  $S_2$ , there exists  $M_0(n, k, \alpha, \epsilon)$  such that

$$E[d_{1,k}^{2\alpha} | X_1 \in S_2] \leq c_1 E[(p(X_1)M)^{-2\alpha/n} | X_1 \in S_2] \leq c_1 \epsilon^{-2\alpha/n} \log^{-\alpha} M \quad (43)$$

for some constant  $c_1(n, k, \alpha)$  and all  $M > M_0$ . We want to bound  $P(d_{1,k} > L\epsilon^{-1/n} \log^{-1/2} M | X_1 \in S_2)$  in order to finish the proof. By Lemma 2, we have for  $0 < r < 1$  and  $x \in S_2$ ,

$$\omega_x(r) \geq c_2 p(x) r^n \geq \frac{c_2 \epsilon r^n \log^{n/2} M}{M} \quad (44)$$

for some constant  $c_2(n)$ . Then because  $\omega_x(r)$  is strictly increasing with respect to  $r$ , using Lemma 1 we have

$$\begin{aligned} P(d_{1,k} > \frac{L}{\epsilon^{1/n} \sqrt{\log M}} | X_1 \in S_2) &= P(\omega_{X_1}(d_{1,k}) > \omega_{X_1}\left(\frac{L}{\epsilon^{1/n} \sqrt{\log M}}\right) | X_1 \in S_2) \\ &\leq P(\omega_{X_1}(d_{1,k}) > \frac{c_2 L^n}{M} | X_1 \in S_2) = k \binom{M-1}{k} \int_{\frac{c_2 L^n}{M}}^1 \omega^{k-1} (1-\omega)^{M-k-1} d\omega \end{aligned}$$

when  $c_2 L^n M^{-1} < 1$  (which can be imposed by taking a sufficiently large threshold  $M_0$ ). We use  $(1-\omega) \leq e^{-\omega}$  and

$$\binom{M-1}{k} \leq \frac{M^k}{k!}$$

to obtain for  $M > c_2 L^n + 4k$ ,

$$k \binom{M-1}{k} \int_{c_2 L^n M^{-1}}^1 \omega^{k-1} (1-\omega)^{M-k-1} d\omega \leq M^k \int_{c_2 L^n M^{-1}}^1 \omega^{k-1} e^{-\frac{1}{2}M\omega} d\omega.$$

The last integral can be bounded by integration by parts as in (19):

$$M^k \int_{c_2 L^n M^{-1}}^1 \omega^{k-1} e^{-\frac{1}{2}M\omega} d\omega \leq 2^k k! (c_2 L^n + 1)^k e^{-\frac{1}{2}c_2 L^n} \leq 4^k k! c_2^k L^{nk} e^{-\frac{1}{2}c_2 L^n}$$

assuming without losing generality that  $c_2 L^n \geq 1$ . We conclude that

$$P(d_{1,k} > \frac{L}{\epsilon^{1/n} \sqrt{\log M}} | X_1 = x) \leq 4^k k! c_2^k L^{nk} e^{-\frac{1}{2}c_2 L^n}. \quad (45)$$



In light of (42), (43) and (45) we have arrived to the conclusion

$$E[d_{1,k}^\alpha (1 - I_L) | X_1 \in S_2] \leq 2^k \sqrt{k! c_1 c_2^k} L^{nk/2} e^{-\alpha/n} e^{-\frac{1}{4}c_2 L^n} \log^{-\alpha/2} M.$$

The term  $L^{nk/2}$  can be dropped in the final conclusion, as it is negligible compared to the exponential decay with respect to  $L$ .  $\square$

The variable  $Y$  emerged in Equation (8). It was defined by

$$Y = \frac{Mp(X_1)}{\log^{n/2} M}. \quad (46)$$

A major idea behind our proofs is the asymptotic uniformity of  $Y$  as shown by

**Lemma 12.** *Suppose that (2) holds. Let  $h(y)$  be a measurable function  $[\epsilon, \epsilon^{-1}] \mapsto [0, 1]$ . Then*

$$E[h(Y) | X_1 \in S_2] \rightarrow \frac{\epsilon}{1 - \epsilon^2} \int_{\epsilon}^{\epsilon^{-1}} h(y) dy.$$

in the limit  $M \rightarrow \infty$ .

*Proof.* The function

$$s(y) = \frac{Me^{-\frac{1}{2}y^2}}{(2\pi)^{n/2} \log^{n/2} M} \quad (47)$$

is strictly decreasing on  $y \in [a, b]$  with  $a$  and  $b$  defined in Equations (34) and (35). It has the inverse  $s^{-1} : [\epsilon, \epsilon^{-1}] \mapsto [a, b]$ :

$$s^{-1}(y) = \sqrt{-2 \log y - n \log \log M + 2 \log M - n \log 2\pi}$$

with the first derivative denoted by  $Ds^{-1}$ . Conditionally on  $X_1 \in S_2$ , the variable  $\|X_1\|$  has the density

$$p_{\|X_1\|}(y) = \frac{nV_n}{(2\pi)^{n/2} P(X_1 \in S_2)} y^{n-1} e^{-\frac{1}{2}y^2}$$

and  $Y$  has the density (on  $[\epsilon, \epsilon^{-1}]$ )

$$p_{\|X_1\|}(s^{-1}(y)) |Ds^{-1}(y)| = \frac{nV_n y s^{-1}(y)^{n-1} \log^{n/2} M}{MP(X_1 \in S_2)} |Ds^{-1}(y)|. \quad (48)$$

Because  $y \in [\epsilon, \epsilon^{-1}]$ , we have in the limit  $M \rightarrow \infty$  with everything else fixed,

$$s^{-1}(y)^{n-1} = (2 \log M)^{n/2-1/2} \left( 1 + O\left(\frac{\log \log M}{\log M}\right) \right) \quad (49)$$

and

$$|Ds^{-1}(y)| = \frac{1}{y \sqrt{2 \log M}} \left( 1 + O\left(\frac{\log \log M}{\log M}\right) \right). \quad (50)$$

By Equations (48)-(50) and Lemma 10 we have

$$P_{\|X_1\|}(s^{-1}(y))|Ds^{-1}(y)| = \frac{\epsilon}{1-\epsilon^2} + O\left(\frac{\log^2 \log M}{\log M}\right). \quad (51)$$

This approximation implies that

$$\begin{aligned} E[h(Y)|X_1 \in S_2] &= \frac{\epsilon}{1-\epsilon^2} \int_{\epsilon}^{\epsilon^{-1}} h(y)dy + O\left(\frac{\log^2 \log M}{\log M}\right) \\ &\rightarrow \frac{\epsilon}{1-\epsilon^2} \int_{\epsilon}^{\epsilon^{-1}} h(y)dy \end{aligned}$$

in the limit  $M \rightarrow \infty$ . □

Next we will find out the asymptotic behavior of  $E[d_{1,k}^\alpha | X_1 \in S_2]$ , which together with the approximation for  $P(X_1 \in S_2)$  takes care of region  $S_2$ . The key to the analysis is Lemma 12. The following represents the nearest neighbor distance in terms of the small ball probability and the variable  $Y$ . We invoke the event  $I_L$  to bound  $d_{1,k}$ ;  $L$  stays fixed in this considerations the idea being the limit  $L \rightarrow \infty$  after taking the limit  $M \rightarrow \infty$ .

**Lemma 13.** *Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then*

$$E[d_{1,k}^\alpha I_L | X_1 \in S_2] = \frac{E[f^{-1}\left(\frac{2^{n/2}M\omega_{X_1}(d_{1,k})}{Y}\right)^\alpha I_L | X_1 \in S_2]}{2^{\alpha/2} \log^{\alpha/2} M} + R_1,$$

where  $Y$  is defined in Equation (46) and

$$|R_1| \leq \frac{c \log \log M}{\log^{\alpha/2+1} M}$$

for some constant  $c(n, \alpha, \epsilon, L)$ .

*Proof.* We first collect a few useful facts. If  $x \in S_2$ , then by Lemma 5

$$\|x\|^n \omega_x(r) = p(x)f(\|x\|r) - p(x)R_1 \quad (52)$$

or equivalently

$$r = \frac{f^{-1}\left(\frac{\|x\|^n \omega_x(r)}{p(x)} + R_1\right)}{\|x\|} \quad (53)$$

with  $0 \leq R_1 \leq r^2 f(\|x\|r)$ .  $x \in S_2$  implies

$$c_1^{-1} \sqrt{\log M} \leq \|x\| \leq c_1 \sqrt{\log M} \quad (54)$$

for some constant  $c_1(n, \epsilon)$ . The indicator function  $I_L$  ensures that we only need to consider

$$0 < r < \frac{L}{\epsilon^{1/n} \sqrt{\log M}}.$$

Then by (54)

$$\|x\|r \leq \frac{Lc_1}{\epsilon^{1/n}}. \quad (55)$$

By a Taylor expansion, for any real number  $\beta \in \mathfrak{R}$  and  $x \in S_2$ ,

$$|\|x\|^\beta - (2 \log M)^{\beta/2}| \leq c_2 \log \log M \log^{\beta/2-1} M \quad (56)$$

for some constant  $c_2(n, \epsilon, \beta)$ . Moreover,  $f$  is an increasing continuous function allowing a bound on  $R_1$ :

$$R_1 \leq r^2 f(\|x\|r) \leq \frac{L^2 f\left(\frac{Lc_1}{\epsilon^{1/n}}\right)}{\epsilon^{2/n} \log M} \leq \frac{c_3}{\log M} \quad (57)$$

for  $c_3 = L^2 \epsilon^{-2/n} f\left(\frac{Lc_1}{\epsilon^{1/n}}\right)$ . Having made the preliminary observations, we are ready for the first step towards completing of the proof. We have for  $x \in S_2$  by Equation (53)

$$E[d_{1,k}^\alpha I_L | X_1 = x] = E \left[ \frac{f^{-1} \left( \frac{\|X_1\|^n \omega_{X_1}(d_{1,k})}{p(X_1)} + R_2 \right)^\alpha}{\|X_1\|^\alpha} I_L | X_1 = x \right]$$

with

$$0 \leq R_2 \leq c_3 \log^{-1} M \quad (58)$$

( $R_2$  is  $R_1$  with  $d_{1,k}$  instead of  $r$  and multiplied by  $I_L$ ). The challenging part is to modify the argument for  $f^{-1}$ . We first tackle the easier task of replacing  $\|X_1\|^\alpha$  with a function of  $M$ . To this end, we observe that

$$E[d_{1,k}^\alpha I_L | X_1 = x] = E \left[ \frac{f^{-1} \left( \frac{\|X_1\|^n \omega_{X_1}(d_{1,k})}{p(X_1)} + R_2 \right)^\alpha}{2^{\alpha/2} \log^{\alpha/2} M} I_L | X_1 = x \right] + R_3 \quad (59)$$

with

$$R_3 = E \left[ f^{-1} \left( \frac{\|X_1\|^n \omega_{X_1}(d_{1,k})}{p(X_1)} + R_2 \right)^\alpha (\|X_1\|^{-\alpha} - 2^{-\alpha/2} \log^{-\alpha/2} M) I_L | X_1 = x \right].$$

By Lemma 5 and Equations (52), (55) and (57) we find a constant  $c_4(n, \epsilon, L)$  such that

$$\frac{\|x\|^n \omega_x(r)}{p(x)} + \frac{c_3}{\log M} \leq f(\|x\|r) + \frac{c_3}{\log M} \leq f\left(\frac{Lc_1}{\epsilon^{1/n}}\right) + \frac{c_3}{\log M} \leq c_4 \quad (60)$$

for  $x \in S_2$  and  $0 < r < L\epsilon^{-1/n} \log^{-1/2} M$ . Using the previous inequality and the fact that  $f^{-1}$  is an increasing function together with Equation (56) allows us to bound

$$|R_3| \leq \frac{c_2(n, \epsilon, -\alpha) f^{-1}(c_4) \log \log M}{\log^{\alpha/2+1} M}. \quad (61)$$

We move to the argument for  $f^{-1}$ . Again, it would be useful to get rid of the norm  $\|x\|^n$ . This is achieved by modifying the argument appearing in (59) (due to conditioning, we may use  $x$  instead of  $X_1$  in the expressions):

$$\frac{\|x\|^n \omega_x(d_{1,k})}{p(x)} I_L = \frac{2^{n/2} \omega_x(d_{1,k}) \log^{n/2} M}{p(x)} I_L + R_4,$$

where by Equation (56) (to bound  $\omega_x(d_{1,k})$ , we use Equations (60) and (54))

$$|R_4| = \frac{|||x||^n - 2^{n/2} \log^{n/2} M| \omega_x(d_{1,k})}{p(x)} I_L \leq \frac{c_5 \log \log M}{\log M} \quad (62)$$

for some constant  $c_5(n, \epsilon, L)$ .

In summary, this far we have shown that

$$\begin{aligned} & E[d_{1,k}^\alpha I_L | X_1 = x] \\ &= E \left[ \frac{f^{-1} \left( \frac{2^{n/2} \omega_{X_1}(d_{1,k}) \log^{n/2} M}{p(X_1)} + R_2 + R_4 \right)^\alpha}{2^{\alpha/2} \log^{\alpha/2} M} I_L | X_1 = x \right] + R_3, \end{aligned} \quad (63)$$

where (58), (61) and (62) bound the three correction terms.

While the correction terms  $R_2$  and  $R_4$  are small, they appear inside the argument for  $f^{-1}$ . The best we can say about their effect is

$$\begin{aligned} & E \left[ \left| f^{-1} \left( \frac{2^{n/2} \omega_{X_1}(d_{1,k}) \log^{n/2} M}{p(X_1)} + R_2 + R_4 \right)^\alpha \right. \right. \\ & \quad \left. \left. - f^{-1} \left( \frac{2^{n/2} \omega_{X_1}(d_{1,k}) \log^{n/2} M}{p(X_1)} \right)^\alpha \right| I_L | X_1 = x \right] \\ & \leq (R_2 + R_4) \sup_{t \in [0, 2^{n/2+1} c_4]} |D(f^{-1}(t)^\alpha)| \end{aligned} \quad (64)$$

assuming without losing generality that  $|R_2 + R_4| \leq c_4$ . So, we need to bound the derivative of the function  $f^{-1}(t)^\alpha$  on bounded intervals. We observe that

$$D(f^{-1}(t)^\alpha) = \frac{\alpha f^{-1}(t)^{\alpha-1}}{Df(f^{-1}(t))}. \quad (65)$$

Furthermore,

$$Df(t) = nt^{n-1} \int_{B(0,1)} e^{ty^{(1)}} dy + t^n \int_{B(0,1)} y^{(1)} e^{ty^{(1)}} dy \geq \frac{1}{2} n V_n t^{n-1}, \quad (66)$$

because

$$\int_{B(0,1)} y^{(1)} e^{ty^{(1)}} dy = \int_0^t \int_{B(0,1)} (y^{(1)})^2 e^{sy^{(1)}} dy ds \geq 0.$$

Using (66) in (65) together with the fact that  $f^{-1}(t)^{\alpha-n}$  ( $\alpha > n$ ) is an increasing function yields

$$\sup_{t \in [0, 2^{n/2+1} c_4]} |D(f^{-1}(t)^\alpha)| \leq \sup_{t \in [0, 2^{n/2+1} c_4]} \frac{2\alpha}{n V_n} f^{-1}(t)^{\alpha-n} \leq \frac{2\alpha f^{-1}(2^{n/2+1} c_4)^{\alpha-n}}{n V_n}.$$

Using the upper bound in (64) shows that for  $x \in S_2$ ,

$$\begin{aligned} & E \left[ \frac{f^{-1} \left( \frac{2^{n/2} \omega_{X_1}(d_{1,k}) \log^{n/2} M}{p(X_1)} + R_2 + R_4 \right)^\alpha}{2^{\alpha/2} \log^{\alpha/2} M} I_L | X_1 = x \right] \\ &= E \left[ \frac{f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha}{2^{\alpha/2} \log^{\alpha/2} M} I_L | X_1 = x \right] + R_5 \end{aligned}$$

with  $|R_5| \leq c_6 \log \log M \log^{-\alpha/2-1} M$  for some constant  $c_6(n, \alpha, \epsilon, L)$ . The proof is finished by recalling the earlier observation (63).  $\square$

In Lemma 13, we find the term  $Y$ , which has the asymptotic uniformity property as proven in Lemma 12. Connecting the two results mainly involves removing the truncation  $I_L$ , but takes some technical effort. The function  $g$  was defined in Equation (25).

**Lemma 14.** *Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then in the limit  $M \rightarrow \infty$*

$$(2 \log M)^{\alpha/2} E[d_{1,k}^\alpha | X_1 \in S_2] \rightarrow \frac{\epsilon}{(k-1)!(1-\epsilon^2)} \int_\epsilon^{\epsilon^{-1}} g\left(\frac{2^{n/2}}{y}\right) dy.$$

*Proof.* By Lemma 13, we know that

$$\begin{aligned} & (2 \log M)^{\alpha/2} E[d_{1,k}^\alpha I_L | X_1 \in S_2] \\ & - E \left[ f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha I_L | X_1 \in S_2 \right] \rightarrow 0 \end{aligned}$$

in the limit  $M \rightarrow \infty$  with  $(n, k, \alpha, \epsilon, L)$  fixed. We write

$$\begin{aligned} & E \left[ f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha I_L | X_1 \in S_2 \right] \\ &= \frac{\int_{S_2} E \left[ f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha I_L | X_1 = x \right] p(x) dx}{P(X_1 \in S_2)}. \end{aligned}$$

Using Equation (47) and Lemma 1 (recall that  $Y$  depends only on  $X_1$ ),

$$\begin{aligned} & E \left[ f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha I_L | X_1 = x \right] \\ &= k \binom{M-1}{k} \int_0^{\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)} \omega^{k-1} (1-\omega)^{M-k-1} f^{-1} \left( \frac{2^{n/2} M \omega}{s(\|x\|)} \right)^\alpha d\omega. \end{aligned}$$

Now

$$k \binom{M-1}{k} = \frac{(M-1)!}{(k-1)!(M-1-k)!} = \frac{M^k}{(k-1)!} + R_1 \quad (67)$$

with  $|R_1| \leq c_1 M^{k-1}$  for some constant  $c_1(k)$ . Also, because  $\|x\|$  behaves asymptotically as  $\sqrt{2 \log M}$  and  $p(x) > \frac{\epsilon \log^{n/2} M}{M}$  on  $S_2$ , Equation (60) shows that

$$\omega_x \left( \frac{L}{\epsilon^{1/n} \sqrt{\log M}} \right) \leq \frac{c_2}{M} \quad (68)$$

for some constant  $c_2(n, \epsilon, L)$ . This implies that for  $\omega < \omega_x(L\epsilon^{-1/n} \log^{-1/2} M)$ ,

$$(1 - \omega)^{M-k-1} = e^{-M\omega} + R_2 \quad (69)$$

with

$$|R_2| \leq |(1 - \omega)^{M-k-1} - e^{-M\omega}| = e^{(M-k-1)\log(1-\omega)+M\omega} - 1 \leq \frac{c_3}{M}$$

for some constant  $c_3(n, k, \alpha, \epsilon, L)$ . By Equations (67)-(69) together with the fact that  $f^{-1}$  is an increasing function,

$$\begin{aligned} & k \binom{M-1}{k} \int_0^{\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)} \omega^{k-1} (1 - \omega)^{M-k-1} f^{-1} \left( \frac{2^{n/2} M \omega}{s(\|x\|)} \right)^\alpha d\omega \\ &= \frac{M^k}{(k-1)!} \int_0^{\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)} \omega^{k-1} e^{-M\omega} f^{-1} \left( \frac{2^{n/2} M \omega}{s(\|x\|)} \right)^\alpha d\omega \\ &+ R_3 + R_4 \end{aligned}$$

with

$$\begin{aligned} |R_3| &\leq \left| k \binom{M-1}{k} - \frac{M^k}{(k-1)!} \right| \int_0^{\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)} \omega^{k-1} \\ &\quad \times (1 - \omega)^{M-k-1} f^{-1} \left( \frac{2^{n/2} M \omega}{\epsilon} \right)^\alpha d\omega \\ &\leq c_1 M^{k-1} \int_0^{c_2 M^{-1}} \omega^{k-1} f^{-1} \left( \frac{2^{n/2} M \omega}{\epsilon} \right)^\alpha d\omega \leq \frac{c_1 c_2^k f^{-1} \left( \frac{2^{n/2} c_2}{\epsilon} \right)^\alpha}{kM} \end{aligned}$$

and

$$\begin{aligned} |R_4| &\leq \frac{M^k}{(k-1)!} \int_0^{c_2 M^{-1}} \omega^{k-1} |e^{-M\omega} - (1 - \omega)^{M-k-1}| f^{-1} \left( \frac{2^{n/2} M \omega}{\epsilon} \right)^\alpha d\omega \\ &\leq \frac{c_2^k c_3 f^{-1} \left( \frac{2^{n/2} c_2}{\epsilon} \right)^\alpha}{M}. \end{aligned}$$

Observe that the bounds for  $R_3$  and  $R_4$  hold for any  $x \in S_2$ . By a change of variables,

$$\begin{aligned} & \frac{M^k}{(k-1)!} \int_0^{\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)} \omega^{k-1} e^{-M\omega} f^{-1} \left( \frac{2^{n/2} M \omega}{s(\|x\|)} \right)^\alpha d\omega \\ &= \frac{1}{(k-1)!} \int_0^\infty \omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{s(\|x\|)} \right)^\alpha d\omega \\ &\quad - \frac{1}{(k-1)!} \int_{M\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)}^\infty \omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{s(\|x\|)} \right)^\alpha d\omega \\ &= \frac{1}{(k-1)!} g \left( \frac{2^{n/2}}{s(\|x\|)} \right) + R_5. \end{aligned}$$

We would like to show that

$$\begin{aligned} & \lim_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} \sup_{x \in S_2} R_5 \\ &= \lim_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} \sup_{x \in S_2} \int_{M\omega_x(L\epsilon^{-1/n} \log^{-1/2} M)}^\infty \omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{s(\|x\|)} \right)^\alpha d\omega = 0. \end{aligned} \quad (70)$$

To see that this is true, we observe that by Lemma 6, for some constant  $c_4(n, k, \alpha, \epsilon)$  there is the bound

$$\omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{s(\|x\|)} \right)^\alpha \leq \omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{\epsilon} \right)^\alpha \leq c_4 \omega^{k-1} (1 + \omega) e^{-\omega}$$

with the upper bound integrable on  $[0, \infty)$  and independent of  $x \in S_2$ . Moreover, by Equation (44)  $\lim_{L \rightarrow \infty} \liminf_{M \rightarrow \infty} \inf_{x \in S_2} M \omega_x \left( \frac{L}{\epsilon^{1/n} \sqrt{\log M}} \right) = \infty$  showing that (70) indeed holds.

In summary, we have shown that

$$\begin{aligned} & \lim_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} E \left[ f^{-1} \left( \frac{2^{n/2} M \omega_{X_1}(d_{1,k})}{Y} \right)^\alpha I_L | X_1 \in S_2 \right] \\ &= \lim_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} \frac{\frac{1}{(k-1)!} \int_{S_2} \int_0^\infty \omega^{k-1} e^{-\omega} f^{-1} \left( \frac{2^{n/2} \omega}{s(\|x\|)} \right)^\alpha p(x) d\omega dx}{P(X_1 \in S_2)} \\ &\quad + \frac{\int_{S_2} (R_3 + R_4 + R_5) p(x) dx}{P(X_1 \in S_2)} = \limsup_{M \rightarrow \infty} \frac{1}{(k-1)!} E \left[ g \left( \frac{2^{n/2}}{Y} \right) | X_1 \in S_2 \right] \end{aligned}$$

and similarly with  $\liminf$  instead of  $\limsup$ . The last limit exists by Lemma 12, which shows that

$$E \left[ g \left( \frac{2^{n/2}}{Y} \right) | X_1 \in S_2 \right] \rightarrow \frac{\epsilon}{1 - \epsilon^2} \int_\epsilon^{\epsilon^{-1}} g \left( \frac{2^{n/2}}{y} \right) dy$$

in the limit  $M \rightarrow \infty$ . On the other hand, Lemma 11 shows that

$$\lim_{L \rightarrow \infty} \limsup_{M \rightarrow \infty} (2 \log M)^{\alpha/2} |E[d_{1,k}^\alpha I_L | X_1 \in S_2] - E[d_{1,k}^\alpha | X_1 \in S_2]| = 0$$

finalizing the proof. □

Now we are able to put everything together to conclude region  $S_2$ :

**Lemma 15.** Assume that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \lim_{M \rightarrow \infty} M \log^{\alpha/2+1-n} M \int_{S_2} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \\ = \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!} \int_0^\infty g\left(\frac{1}{y}\right) dy < \infty. \end{aligned}$$

*Proof.* The claim follows from Lemmas 10 and 14:

$$\begin{aligned} M \log^{\alpha/2+1-n} M \int_{S_2} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \\ \rightarrow \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!(1-\epsilon^2)} \int_{2^{-n/2}\epsilon}^{2^{-n/2}\epsilon^{-1}} g\left(\frac{1}{y}\right) dy \end{aligned}$$

in the limit  $M \rightarrow \infty$ . To finish the proof, we would like to replace the integration limits by 0 and  $\infty$  when  $\epsilon \rightarrow 0$ , which amounts to showing that  $g(y^{-1})$  is an integrable function. Integrability can be established using Lemma 6 to show that

$$\int_0^\infty g\left(\frac{1}{y}\right) dy \leq c \int_0^1 (1 + \log^\alpha y^{-1}) dy + c \int_1^\infty y^{-\alpha/n} dy$$

for some constant  $c(n, k, \alpha)$ . Both terms in the right side are finite (the second one because  $\alpha > n$ ).  $\square$

## 8 Region $S_3$

$S_3$  consists of points where the density  $p$  takes small values:

$$S_3 = \{x \in \mathfrak{R}^n : p(x) < \frac{\epsilon \log^{n/2} M}{M}\}.$$

To bound nearest neighbor distances on  $S_3$  we need similar tools as for  $S_2$ , but only upper bounds are needed providing some more flexibility. The sets  $\tilde{S}_{3,i}$  are defined analogously to (27):

$$\tilde{S}_{3,i} = \{x \in \mathfrak{R}^n : 2^{-i-1} \frac{\epsilon \log^{n/2} M}{M} \leq p(x) < 2^{-i} \frac{\epsilon \log^{n/2} M}{M}\}$$

for  $0 \leq i \leq i^*$  with

$$i^* = \left\lceil \frac{(\alpha + 1)}{\log 2} \log \log M \right\rceil + 1.$$

Moreover,  $\tilde{S}_{3,C} = S_3 \setminus \cup_{i=0}^{i^*} \tilde{S}_{3,i}$ . Then we have



**Lemma 16.** Assume that  $p(x)$  is the multivariate Gaussian distribution (2). Then for some threshold  $M_0(n, \epsilon)$  we have for  $M > M_0$  and  $0 \leq i \leq i^*$  that

$$P(X_1 \in \tilde{S}_{3,i}) \leq 2^{-i} c \frac{\epsilon \log^{n-1} M}{M}$$

for some constant  $c(n)$ .

*Proof.* The set  $\tilde{S}_{3,i}$  consists of points  $x \in \mathfrak{R}^n$  with  $\|x\| \in [a, b]$  and

$$\begin{aligned} a &= \sqrt{2 \log M - n \log \log M - 2 \log \epsilon + i \log 4 - n \log(2\pi)} \\ b &= \sqrt{2 \log M - n \log \log M - 2 \log \epsilon + (i+1) \log 4 - n \log(2\pi)}. \end{aligned} \quad (71)$$

Using the mean value theorem for  $a$  and  $b$  we have for  $0 \leq i \leq i^*$ ,

$$b - a \leq \frac{4}{\sqrt{\log M}} \quad (72)$$

after some threshold  $M_0(n, \epsilon)$ . Also, we may take  $b \leq \sqrt{3 \log M}$  for  $0 \leq i \leq i^*$  as the term  $2 \log M$  inside the square root (71) grows faster than the other terms. Then

$$\lambda(\tilde{S}_{3,i}) = nV_n \int_a^b x^{n-1} dx \leq nV_n b^{n-1} (b-a) \leq 3^{n/2+3/2} nV_n \log^{n/2-1} M. \quad (73)$$

By Equation (73) and the fact  $p(x) \leq 2^{-i} \epsilon M^{-1} \log^{n/2} M$  on  $\tilde{S}_{3,i}$ , we have

$$P(X_1 \in \tilde{S}_{3,i}) \leq 2^{-i} 3^{n/2+3/2} nV_n \frac{\epsilon \log^{n-1} M}{M}.$$

□

Assessing the contributions from  $\tilde{S}_{3,i}$  is convenient by using the function  $f$  together with the small ball probability. The proof idea is essentially similar to that used for  $S_2$  in Section 7, but because we need only an upper bound the proof is easier.

**Lemma 17.** Suppose that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then for some threshold  $M_0(n, \alpha, k, \epsilon)$ , we have for  $M > M_0$  and  $0 \leq i \leq i^*$  that

$$\int_{\tilde{S}_{3,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c 2^{-i} \epsilon (\log^\alpha \epsilon^{-1} + i^\alpha + 1) \frac{\log^{n-\alpha/2-1} M}{M}$$

for some constant  $c(n, k, \alpha)$ .

*Proof.* We decompose

$$\int_{\tilde{S}_{3,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx = (I_1 + I_2) P(X_1 \in \tilde{S}_{3,i})$$

with

$$I_1 = E[d_{1,k}^\alpha I(d_{1,k} \leq 1) | X_1 \in \tilde{S}_{3,i}]$$

$$I_2 = E[d_{1,k}^\alpha I(d_{1,k} > 1) | X_1 \in \tilde{S}_{3,i}].$$

$P(X_1 \in \tilde{S}_{3,i})$  was computed in Lemma 16.

### 1. The term $I_1$

If  $0 < r < 1$ , we have

$$\begin{aligned} \|x\|^n \omega_x(r) &= (2\pi)^{-n/2} \|x\|^n \int_{B(x,r)} e^{-\frac{1}{2}\|y\|^2} dy \\ &\geq e^{-\frac{1}{2}\|x\|^n} p(x) \int_{B(0,r)} e^{-x^T y} dy = e^{-\frac{1}{2}\|x\|^n} p(x) f(\|x\|r), \end{aligned} \quad (74)$$

where the function  $f$  was defined in Lemma 5. This implies that

$$d_{1,k} \leq \frac{f^{-1}\left(\frac{e^{\frac{1}{2}\|X_1\|^n \omega_{X_1}(d_{1,k})}}{p(X_1)}\right)}{\|X_1\|}. \quad (75)$$

By taking  $M_0$  large enough, we may ensure that

$$\sqrt{\log M} \leq \|x\| \leq \sqrt{3 \log M} \quad (76)$$

on  $x \in \tilde{S}_{3,i}$  for  $0 \leq i \leq i^*$ . Then by Lemma 6 and Equations (75)-(76),

$$E[d_{1,k}^\alpha I(d_{1,k} \leq 1) | X_1 = x] \leq c_1 E \left[ \frac{1 + \log^\alpha \left( 1 + \frac{2^{i+n+2} M \omega_{X_1}(d_{1,k})}{\epsilon} \right)}{\log^{\alpha/2} M} \middle| X_1 = x \right]$$

for some constant  $c_1(n, \alpha)$ . Using  $\log(1+z) \leq z$  for  $z \geq 0$ , we have

$$\log(1 + 2^{i+n+2} \epsilon^{-1} M \omega_{X_1}(d_{1,k})) \leq (i+n+2) \log 2 + \log \epsilon^{-1} + M \omega_{X_1}(d_{1,k})$$

recalling that  $0 < \epsilon < 1$ . The  $\alpha$ -moment of the conditional expectation of the last expression is bounded by  $c_2(\log^\alpha \epsilon^{-1} + i^\alpha + 1)$  for some constant  $c_2(n, k, \alpha)$  by Lemma 1 and Equation (13).

### 2. The term $I_2$

By Hölder's inequality, Lemma 3 and Equation (76),

$$I_2 \leq c_3 \log^{\alpha/2} M \sqrt{P(d_{1,k} > 1 | X_1 \in \tilde{S}_{3,i})}$$

for some constant  $c_3(n, k, \alpha)$ . Equation (19) can be applied here: for  $x \in \tilde{S}_{3,i}$ ,

$$P(d_{1,k} > 1 | X_1 = x) \leq c_4(n, k) e^{-c_4(n,k)^{-1} M \omega_x(1)}$$

for some constant  $c_4(n, k)$ . It would be sufficient to show that for any  $j > 0$ ,

$$\sup_{0 \leq i \leq i^*, x \in \tilde{S}_{3,i}} \omega_x(1) M \log^{-j} M \rightarrow \infty \quad (77)$$

in the limit  $M \rightarrow \infty$ . By Equations (74) and (76) taking into account that on  $\tilde{S}_{3,i}$ ,

$$p(x) \geq 2^{-i^*} \frac{\epsilon \log^{n/2} M}{M} \geq \frac{\epsilon \log^{n/2-\alpha-1} M}{2M},$$

we have

$$\begin{aligned} \omega_x(1) &\geq e^{-\frac{1}{2}} \frac{\epsilon}{4M \log^{\alpha+1-n/2} M} \int_{B(0,1)} e^{\sqrt{\log M} y^{(1)}} dy \\ &\geq e^{-\frac{1}{2}} \frac{\epsilon}{4M \log^{\alpha+1-n/2} M} \lambda(B(0,1) \cap \{y \in \mathfrak{R}^n : y^{(1)} > \frac{1}{2}\}) e^{\frac{1}{2} \sqrt{\log M}}. \end{aligned}$$

The term  $e^{\frac{1}{2} \sqrt{\log M}}$  approaches infinity faster than  $\log^j M$  for any  $j > 0$ . This shows (77) and we conclude that for any  $j > 0$ ,  $I_2$  approaches 0 faster than  $\log^{-j} M$  in the limit  $M \rightarrow \infty$ .  $\square$

The region  $\tilde{S}_{3,C}$  is easier, because by taking  $i^*$  as a large number, we are able to control the probability of this set.

**Lemma 18.** *Suppose that  $p(x)$  is the multivariate Gaussian distribution (2) and  $\alpha > n$ . Then there exists a constant  $c(n, k, \alpha)$  and a threshold  $M_0(n, k, \alpha, \epsilon)$  such that for  $M > M_0$ , we have*

$$\int_{\tilde{S}_{3,C}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c \epsilon \frac{\log^{n-\alpha/2-1} M}{M}.$$

*Proof.* On  $\tilde{S}_{3,C}$  we have

$$p(x) \leq 2^{-i^*} \frac{\epsilon \log^{n/2} M}{M} \leq \frac{\epsilon \log^{n/2-\alpha-1} M}{M}.$$

We define

$$T_i = 2^{i+1} \tilde{S}_{3,C} \setminus 2^i \tilde{S}_{3,C},$$

where  $2^i \tilde{S}_{3,C} = \{x \in \mathfrak{R}^n : 2^{-i} x \in \tilde{S}_{3,C}\}$ . We may assume that  $\|x\| \leq 4\sqrt{\log M}$  on  $T_0$  and consequently  $\|x\| \leq 2^{i+2} \sqrt{\log M}$  on  $T_i$  for any  $i \geq 0$ .  $\lambda(T_i)$  is roughly bounded by  $\lambda(T_i) \leq 2^{n(i+2)} V_n \log^{n/2} M$ . Now by Lemma 3,

$$\begin{aligned} \int_{\tilde{S}_{3,C}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx &\leq c_1 \sum_{i=0}^{\infty} \int_{T_i} (\|x\|^\alpha + 1) p(x) dx \\ &\leq \sum_{i=0}^{\infty} 2^{i\alpha} c_2 \log^{\alpha/2} M \left( \frac{(2\pi)^{n/2} \epsilon \log^{n/2-\alpha-1} M}{M} \right)^{2^{2i}} \lambda(T_i) \\ &\leq c_2 V_n \epsilon \frac{\log^{n-\alpha/2-1} M}{M} \sum_{i=0}^{\infty} 2^{(n+\alpha)i+2n} \left( \frac{(2\pi)^{n/2} \epsilon \log^{n/2-\alpha-1} M}{M} \right)^{2^{2i}-1}, \end{aligned} \quad (78)$$

where  $c_1(n)$  is some constant, and to be exact,  $c_2 = 2(2\pi)^{-n/2}c_1$ . The factor 2 in  $c_2$  comes from the fact that  $\log^{\alpha/2} M > 1$  for  $M > 3$  (which can be assumed without losing generality). Now it is rather obvious that the sum does not pose problems.  $\square$

**Lemma 19.** *Assume that (2) holds,  $\alpha > n$  and  $\epsilon < 1/2$  (only small values of  $\epsilon$  matter in any case). Then there exists a threshold  $M_0(n, k, \alpha, \epsilon)$  such that for any  $M > M_0(n, k, \alpha, \epsilon)$ , we have*

$$\int_{S_3} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c\epsilon \log \epsilon^{-1} \frac{\log^{n-\alpha/2-1} M}{M}$$

for some constant  $c(n, k, \alpha)$ .

*Proof.* We decompose

$$\begin{aligned} \int_{S_3} E[d_{1,k}^\alpha | X_1 = x] p(x) dx &= \sum_{i=0}^{i^*} \int_{\tilde{S}_{3,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \\ &\quad + \int_{\tilde{S}_{3,c}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx. \end{aligned}$$

By Lemma 17,

$$\sum_{i=0}^{i^*} \int_{\tilde{S}_{3,i}} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \leq c\epsilon \frac{\log^{n-\alpha/2-1} M}{M} \sum_{i=0}^{\infty} 2^{-i} (\log^\alpha \epsilon^{-1} + i^\alpha + 1).$$

Lemma 18 finalizes the proof.  $\square$

## 9 Proof of Theorem 2

Previously we have examined the regions  $S_1$ ,  $S_2$  and  $S_3$ , which were defined in terms of  $\epsilon$  and  $M$ . We decompose

$$(M \log^{\alpha/2+1-n} M) E[d_{1,k}^\alpha] = I_{1,\epsilon,M} + I_{2,\epsilon,M} + I_{3,\epsilon,M}$$

with

$$I_{i,\epsilon,M} = M \log^{\alpha/2+1-n} M \int_{S_i} E[d_{1,k}^\alpha | X_1 = x] p(x) dx \quad (i = 1, 2, 3).$$

Lemmas 9 and 19 show that  $\lim_{\epsilon \rightarrow 0} \limsup_{M \rightarrow \infty} I_{1,\epsilon,M} + I_{3,\epsilon,M} = 0$ . Also by Lemma 15,

$$\lim_{\epsilon \rightarrow 0} \lim_{M \rightarrow \infty} I_{2,\epsilon,M} = \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!} \int_0^\infty g\left(\frac{1}{y}\right) dy.$$

We conclude that

$$\lim_{M \rightarrow \infty} (M \log^{\alpha/2+1-n} M) E[d_{1,k}^\alpha] = \frac{2^{n-\alpha/2-1} n V_n}{(k-1)!} \int_0^\infty g\left(\frac{1}{y}\right) dy.$$

## References

- [1] D. Evans and A. J. Jones, *A proof of the Gamma test*, Proceedings of the Royal Society A **458** (2002), no. 2027, 2759–2799. MR1942807
- [2] D. Evans, A. J. Jones, and W. M. Schmidt, *Asymptotic moments of near-neighbour distance distributions*, Proceedings of the Royal Society A **458** (2002), no. 2028, 2839–2849. MR1987515
- [3] J. Hansen and I. R. McDonald, *Theory of simple liquids*, Academic Press, 1990.
- [4] M. Kohler, A. Krzyzak, and H. Walk, *Rates of convergence for partitioning and nearest neighbor regression estimates with unbounded data*, Journal of Multivariate Analysis **97** (2006), no. 2, 311–323. MR2234025
- [5] L. F. Kozachenko and N. N. Leonenko, *Sample estimate of entropy of a random vector*, Problems of Information Transmission **23** (1987), no. 2, 9–16. MR0908626
- [6] N. Leonenko and L. Pronzato, *Correction: A class of Rényi information estimators for multidimensional densities*, Annals of Statistics **38** (2010), no. 6, 3837–3838. MR2766870
- [7] N. N. Leonenko, L. Pronzato, and V. Savani, *A class of Rényi information estimators for multidimensional densities*, Annals of Statistics **36** (2008), no. 5, 2153–2182. MR2458183
- [8] E. Liitiäinen, A. Lendasse, and F. Corona, *A boundary corrected expansion of the moments of nearest neighbor distributions*, Random Structures and Algorithms **37** (2010), no. 2, 223–247. MR2676030
- [9] E. Liitiäinen, Amaury Lendasse, and Francesco Corona, *On the statistical estimation of Rényi entropies*, IEEE Conference on Machine Learning for Signal Processing, 2009.
- [10] M. D. Penrose, *Random geometric graphs*, Oxford Studies in Probability, no. 5, Oxford University Press, 2003. MR1986198
- [11] ———, *Laws of large numbers in stochastic geometry with statistical applications*, Bernoulli **13** (2007), no. 4, 1124–1150. MR2364229
- [12] M. D. Penrose and J. E. Yukich, *Laws of large numbers and nearest neighbor distances*, Advances in Directional and Linear Statistics (M. T. Wells and A. SenGupta, eds.), Physica-Verlag HD, 2011, pp. 189–199. MR2767541
- [13] B. Ranneby, S. R. Jammalamadaka, and A. Teterukovskiy, *The maximum spacing estimation for multivariate observations*, Journal of Statistical Planning and Inference **129** (2005), no. 1-2, 427–446. MR2126858
- [14] A. R. Wade, *Explicit laws of large numbers for random nearest-neighbour type graphs*, Advances in Applied Probability **39** (2007), no. 2, 326–342. MR2341576
- [15] J. Yukich, *Probability theory of classical Euclidean optimization problems*, Lecture Notes in Mathematics, Springer, 1998. MR1632875