



# Random Dimension Low Sample Size Asymptotics

Gerd Christoph<sup>1</sup> and Vladimir V. Ulyanov<sup>2,3</sup>(✉)

<sup>1</sup> Department of Mathematics, University of Magdeburg, Magdeburg, Germany  
gerd.christoph@ovgu.de

<sup>2</sup> Faculty of Computational Mathematics and Cybernetics,  
Lomonosov Moscow State University, Moscow, Russian Federation  
vulyanov@cs.msu.ru

<sup>3</sup> Faculty of Computer Science, HSE University,  
109028 Moscow, Russian Federation

**Abstract.** A first investigation of high-dimensional low-sample-size (HDLSS) asymptotics, Hall, Marron and Neeman (2005) discovered a surprisingly rigid geometric structure. A sample of size  $k$  taken from the standard  $m$ -dimensional normal distribution is for large  $m$  close to the vertices of the  $k$ -dimensional simplex in  $m$ -dimensional vector space. It follows from the analysis of three geometric statistics: the length of an observation, the distance between any two independent observations and the angle between these vectors. We generalize and refine the results constructing the second order Chebyshev-Edgeworth expansions under assumption that the data dimension is random and different scaling factors are chosen.

**Keywords:** HDLSS data · Chebyshev-Edgeworth expansions · Random dimension · Student's  $t$ -distribution · Laplace approximation

## 1 Three Geometric Statistics of Gaussian Vectors

We continue to study properties of high-dimensional Gaussian random vectors. In our earlier papers Christoph, Prokhorov and Ulyanov [8] and Bobkov, Naumov and Ulyanov [5] two-sided bounds were constructed for a probability density function of the distance of a Gaussian random element  $Y$  with zero mean from a point  $a$  in a Hilbert space  $\mathbb{H}$ . We get new results for basic geometric statistics connected with high-dimensional random normal vectors.

Let  $\mathbf{X}_1 = (X_{1,1}, \dots, X_{1,m})^T, \dots, \mathbf{X}_k = (X_{k,1}, \dots, X_{k,m})^T$  be a random sample.

In a high-dimension low-sample-size (HDLSS) data it is assumed that dimension  $m$  tends to infinity and sample size  $k$  is fixed.

One of the first investigation of HDLSS data was done in Hall, Marron and Neeman (2005) [14]. It became the basis of research in high-dimensional mathematical statistics. See a recent survey on HDLSS asymptotics and its applications in Aoshima et al. [1]. Further development see e.g. in Fujikoshi, Ulyanov

and Shimizu [12] when both  $m$  and  $k$  may tend to infinity. This is an important framework of the current data analysis called *Big data*. In [14] it was discovered a surprisingly rigid geometric structure. A sample of size  $k$  taken from the standard  $m$ -dimensional normal distribution is close for large  $m$  to the vertices of the  $k$ -dimensional simplex in  $\mathbb{R}^m$ . It follows from the analysis of three geometric statistics:

the **length**  $\|\mathbf{X}_i\|_m$  of an observation,

the **distance**  $\|\mathbf{X}_i - \mathbf{X}_j\|_m$  between any two independent observations,  
and the **angle**  $\theta_m = \text{ang}(\mathbf{X}_i, \mathbf{X}_j)$  between these vectors.

We generalize and refine the results constructing the second order Chebyshev-Edgeworth expansions under assumption that the data dimension is random and different scaling factors are chosen.

In case of  $\dim \mathbb{H} < \infty$  we consider a sample of size  $k$  when the dimension of the observations is a random variable  $N_n$  with values in  $\mathbb{N}_+ = \{1, 2, \dots\}$ .

The present work continues our investigations in Christoph and Ulyanov [9] on these three geometric statistics of Gaussian vectors with randomly distributed dimension  $N_n$  which depends on parameter  $n \in \mathbb{N}_+$  and  $N_n \rightarrow \infty$  in probability as  $n \rightarrow \infty$ . Let the vectors  $\mathbf{X}_1, \dots, \mathbf{X}_k$  and  $N_1, N_2, \dots$  be defined on one and the same probability space and it is assumed that they are independent. If  $T_m := T_m(\mathbf{X}_1, \dots, \mathbf{X}_k)$  is some statistic of the vectors  $\mathbf{X}_1, \dots, \mathbf{X}_k$  with *non-random dimension*  $m \in \mathbb{N}_+$  then the random variable  $T_{N_n} = T_{N_n}(\omega)$  is defined as:

$$T_{N_n}(\omega) := T_{N_n(\omega)}(\mathbf{X}_1(\omega), \dots, \mathbf{X}_k(\omega)), \quad \omega \in \Omega \quad \text{and} \quad n \in \mathbb{N}_+.$$

Therefore, the statistics  $T_{N_n}$  based on statistics  $T_m$  are constructed from the sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_k\}$ , where these vectors have the dimension  $N_n$ .

In [9], the distribution function of the normalized angle  $\theta_m = \text{ang}(\mathbf{X}_i, \mathbf{X}_j)$  was approximated by a second order Chebyshev-Edgeworth expansion with a bound  $\leq Cm^{-2}$  for all  $m \in \mathbb{N}_+$ . Furthermore, the fixed dimension  $m$  of the Gaussian vectors was substituted by a random number  $N_n$  and expansions for statistics  $\theta_{N_n}$  were proved.

A natural question arises whether similar results hold for the length  $\|\mathbf{X}_i\|_{N_n}$  and the distance  $\|\mathbf{X}_i - \mathbf{X}_j\|_{N_n}$  of Gaussian vectors with random dimension  $N_n$ .

Two cases of random dimensions (or random sample sizes)  $N_n$  are considered as e.g. in Bening, Galieva and Korolev [2], Christoph, Monakhov and Ulyanov [7] and Christoph and Ulyanov [9]:

- i) The random dimension  $N_n = N_n(r) \in \mathbb{N}_+$  has negative binomial distribution displaced by 1 with probability of success  $1/n$ , positive parameter  $r > 0$  and probabilities

$$\mathbb{P}(N_n(r) = j) = \frac{\Gamma(j+r-1)}{\Gamma(j)\Gamma(r)} \left(\frac{1}{n}\right)^r \left(1 - \frac{1}{n}\right)^{j-1}, \quad j \in \mathbb{N}_+. \quad (1)$$

- ii) The random dimension  $N_n = N_n(s) \in \mathbb{N}_+$  is discrete Pareto-like distributed with parameters  $n \in \mathbb{N}_+$ ,  $s > 0$  and distribution function

$$\mathbb{P}(N_n(s) \leq k) = \left( \frac{k}{s+k} \right)^n \quad \text{where} \quad N_n(s) = \max_{1 \leq j \leq n} Y_j(s), \quad (2)$$

and  $Y(s), Y_1(s), Y_2(s), \dots$ , are independent discrete Pareto II distributed random variables with the common distribution

$$\mathbb{P}(Y(s) \leq k) = \frac{k}{s+k} \quad \text{and} \quad \mathbb{P}(Y(s) = k) = \frac{s}{(s+k)(s+k-1)}, \quad k \in \mathbb{N}_+. \quad (3)$$

The discrete  $Y(s)$  on integers is the discretized continuous Pareto II (Lomax) random variable, see Buddana and Kozubowski [6].

Both cases of random dimensions of the Gaussian vectors are also interesting because  $\mathbb{E}N_n(r) = r(n-1) + 1 < \infty$  and  $\mathbb{E}N_n(s) = \infty$ , which has an influence on the normalization factors.

The rest of the paper is organized as follows: In Sect. 2, Chebyshev-Edgeworth expansions are proved for the geometric statistics of Gaussian vectors with fixed dimension  $m$ . Section 3 presents the transfer theorem for results with fixed sample size (in our case the dimension of the vectors)  $m$  to those with random sample size  $N_n$ . The main results are given in Sects. 4 and 5 when the random sample size is negative binomial  $N_n(r)$  or discrete Pareto-like  $N_n(s)$  distributed, respectively. In Sect. 6 the main results are proved.

## 2 Approximation for Geometric Statistics of $m$ -Dimensional Normal Vectors

Let  $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,m})^T, \dots, \mathbf{X}_j = (X_{j,1}, \dots, X_{j,m})^T$  be  $m$ -dimensional vectors chosen from a sample  $\{\mathbf{X}_1, \dots, \mathbf{X}_k\}$  of normal distribution  $\mathcal{N}(\mathbf{0}_m, \mathbf{I}_m)$  with mean vectors  $\mathbb{E}\mathbf{X}_k = \mathbf{0}_m$  and covariance matrix  $\mathbf{I}_m$  for  $1 \leq i < j \leq k \leq m$ .

The **length** of the vector  $\mathbf{X}_j$  is defined by the Euclidean distance  $\|\cdot\|_m$ :

$$\|\mathbf{X}_i\|_m = S_m^{1/2} \quad \text{with} \quad S_m = \sum_{k=1}^m X_{i,k}^2. \quad (4)$$

and similarly the **distance**  $\|\mathbf{X}_i - \mathbf{X}_j\|_m$  between any two independent vectors

$$\|\mathbf{X}_i - \mathbf{X}_j\|_m = \sum_{k=1}^m (X_{i,k} - X_{j,k})^2. \quad (5)$$

The distribution of distance  $\|\mathbf{X}_i - \mathbf{X}_j\|_m$  is closely linked to the distribution of length  $\|\mathbf{X}_i\|_m$ , since  $(X_{i,k} - X_{j,k})/\sqrt{2}$  has also standard normal distribution  $\Phi(x)$ . Therefore

$$\mathbb{P}(\|\mathbf{X}_i - \mathbf{X}_j\|_m/\sqrt{2} \leq x) = \mathbb{P}(\|\mathbf{X}_i\|_m \leq x). \quad (6)$$

The **angle**  $\theta_m = \text{ang}(\mathbf{X}_i, \mathbf{X}_j)$  between these two independent vectors with vertex at the origin and the **sample correlation coefficient**  $R_m(\mathbf{X}_i, \mathbf{X}_j)$  are connected by:

$$\cos \theta_m = \frac{\|\mathbf{X}_i\|_m^2 + \|\mathbf{X}_j\|_m^2 - \|\mathbf{X}_i - \mathbf{X}_j\|_m^2}{2\|\mathbf{X}_i\|_m \|\mathbf{X}_j\|_m} = R_m(\mathbf{X}_i, \mathbf{X}_j) = R_m. \quad (7)$$

Hall, Marron and Neeman [14] showed

- for the length  $\|\mathbf{X}_i\|_m = \sqrt{m} + \mathcal{O}_p(1)$ ,
- for the distance  $\|\mathbf{X}_i - \mathbf{X}_j\|_m = \sqrt{2m} + \mathcal{O}_p(1)$  with  $i \neq j$  and
- for the  $\theta_m = \text{angle} \text{ang}(\mathbf{X}_i, \mathbf{X}_j) = \frac{1}{2}\pi + \mathcal{O}_p(m^{-1/2})$  with  $i \neq j$ ,

where  $1 \leq i < j \leq k \leq m$  and  $\mathcal{O}_p$  refers to the stochastic boundedness.

The length of the vector  $\mathbf{X}_i$  drawn from an  $m$ -dimensional normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$  is defined in (4) as  $\|\mathbf{X}_i\|_m = S_m^{1/2}$ , where the statistics  $S_m$  as a sum of the squares of  $m$  independent standard normal random variables has **chi-square distribution** with  $m$  degrees of freedom and

$$V_m = \frac{S_m - m}{\sqrt{2m}} \quad (8)$$

is asymptotically standard normally distributed. With the two-term Chebyshev-Edgeworth expansions in the central limit theorem for the distribution function of  $V_m$ , the following inequality results for all  $m \in \mathbb{N}$

$$\left| P(V_m \leq x) - \Phi(x) - \varphi(x) \left( \frac{\lambda_3 H_2(x)}{6\sqrt{m}} + \frac{\lambda_3^2 H_5(x)}{72m} + \frac{\lambda_4 H_3(x)}{24m} \right) \right| \leq \frac{C}{m^{3/2}}$$

where  $H_2(x) = x^2 - 1$ ,  $H_3(x) = x^3 - 3x$ ,  $H_5(x) = x^5 - 10x^3 + 15x$  are the Chebyshev-Hermite polynomials, skewness  $\lambda_3 = \sqrt{8}$  and excess kurtosis  $\lambda_4 = 12$  of  $S_1$ , see Petrov [19, Sec. 5.7, Theorem 5.18].

Then  $S_m = m(1 + \sqrt{2/m} V_m)$  and Taylor expansion of  $(1 + u)^{1/2}$  lead to

$$\|\mathbf{X}_i\|_m = S_m^{1/2} = \sqrt{m} \left( 1 + \frac{1}{\sqrt{2m}} V_m - \frac{1}{4m} V_m^2 + \frac{\sqrt{2}}{8m^{3/2}} V_m^3 + \dots \right) \quad (9)$$

Define the statistics

$$Z_m = \sqrt{2} \left( \frac{\|\mathbf{X}_i\|_m}{\sqrt{m}} - 1 \right) \quad \text{and} \quad Z_m^* = \sqrt{2} \left( \frac{\|\mathbf{X}_i - \mathbf{X}_j\|_m}{\sqrt{2m}} - 1 \right), \quad (10)$$

then (6) results in

$$P(\sqrt{m} Z_m \leq x) = P(\sqrt{m} Z_m^* \leq x). \quad (11)$$

It follows from (9) that the statistic  $T_1 = \sqrt{m} Z_m$  holds

$$T_1 = \sqrt{m} Z_m = V_m - \frac{\sqrt{2}}{4\sqrt{m}} V_m^2 + \frac{\sqrt{1}}{4m} V_m^3 + \dots \quad (12)$$

Following the sketch of the proof in Kawaguchi, Ulyanov and Fujikoshi [16, Theorem 1] (The coefficients in the polynomial  $l_2(x)$  are incorrect.) and calculating the characteristic function  $f_{T_1}(t)$ , we obtain

$$\begin{aligned}
 f_{T_1}(t) &= \mathbb{E} \left[ e^{itV_m} \left( 1 - \frac{\sqrt{2}(it)}{4\sqrt{m}} V_m^2 + \frac{(it)}{4m} V_m^3 + \frac{(it)^2}{16m} V_m^4 + \mathcal{O}_p(m^{-3/2}) \right) \right] \\
 &= e^{-t^2/2} \left( 1 - \frac{\sqrt{2}((it)^3 + 3(it))}{12\sqrt{m}} + \frac{(it)^6 - 6(it)^4 - 9(it)^2}{144m} \right) + \mathcal{O}(m^{-3/2}). \tag{13}
 \end{aligned}$$

This results in the related expansion of the corresponding distribution function:

**Proposition 1.** *Let  $\mathbf{X}_i$  be a vector drawn from an  $m$ -dimensional normal distribution  $\mathcal{N}(\mathbf{0}_m, \mathbf{I}_m)$ . Then with the asymptotic expansion for the distribution of normalized length  $Z_m = \sqrt{2} \left( \frac{\|\mathbf{X}_i\|_m}{\sqrt{m}} - 1 \right)$  we obtain the following inequality for all  $m \in \mathbb{N}$ :*

$$\left| P\left(\sqrt{m} Z_m \leq x\right) - \Phi(x) - \varphi(x) \left( \frac{x^2 - 4}{6\sqrt{2}m} + \frac{x^5 - 16x^3 + 24x}{144m} \right) \right| \leq \frac{C}{m^{3/2}}. \tag{14}$$

**Corollary 1.** *Let  $\mathbf{X}_i$  and  $\mathbf{X}_j$ ,  $i \neq j$  be independent random vectors with an  $m$ -dimensional normal distribution  $\mathcal{N}(\mathbf{0}_m, \mathbf{I}_m)$ . Due to (11), distribution function of the normalized distance  $Z_m^* = \sqrt{2} \left( \frac{\|\mathbf{X}_i - \mathbf{X}_j\|_m}{\sqrt{2}m} - 1 \right)$  has the same asymptotic expansion as the distribution of normalized length  $Z_m$  and inequality (14) with replacing  $Z_m$  by  $Z_m^*$ .*

Second order Chebyshev-Edgeworth expansion of the angle  $\theta_m = \text{ang}(\mathbf{X}_i, \mathbf{X}_j)$  between independent vectors  $\mathbf{X}_i$  and  $\mathbf{X}_j$  with vertex at the origin and the corresponding sample correlation coefficient  $R_m(\mathbf{X}_i, \mathbf{X}_j)$  with computable error bounds of approximation are shown in Christoph and Ulyanov [9, Section 2], using results of Konishi [17, Sect. 4], Johnson, Kotz and Balakrishnan [15, Chap. 32], Christoph, Ulyanov and Fujikoshi [11]:

$$\sup_x \left| P\left(\sqrt{m} R_m \leq x\right) - \Phi(x) - \frac{x^3 - 5x}{4m} \varphi(x) \right| \leq \frac{B_1}{m^2} \tag{15}$$

and

$$\sup_x \left| P\left(\sqrt{m}(\theta_m - \frac{\pi}{2}) \leq x\right) - \Phi(x) - \frac{x^3 - 15x}{12m} \varphi(x) \right| \leq \frac{B_2}{m^2}. \tag{16}$$

The estimates (15) and (16) were used in Christoph and Ulyanov [9] to obtain second order approximations the statistics  $R_{N_n}$  and  $\Theta_{N_n} = \theta_{N_n} - \pi/2$  when the non-random dimension  $m$  of the vectors is replaced be a random dimension  $N_n$ , where the random dimension  $N_n \rightarrow \infty$  in probability when the parameter  $n \rightarrow \infty$ .

Analogous results for the statistics  $\|\mathbf{X}_i\|_m$  and  $\|\mathbf{X}_i - \mathbf{X}_j\|_m$  are proven in Sects. 4 and 5 below, when the non-random dimension  $m$  is replaced be a random dimension  $N_n$ .

### 3 Auxiliary Proposition

In this section, expansions for the distribution function of statistics  $T_{N_n}$  obtained from samples with random sample size (here with random dimension  $N_n$  of the considered vectors  $\mathbf{X}_i$ ) are obtained. These depend directly on the expansions concerning statistics  $T_m$  based on non-random samples size  $m$  and expansions regarding the random sample size  $N_n$ .

First we formulate the conditions determining expansions for the statistic  $T_m$  with  $\mathbb{E}T_m = 0$  and the normalized random dimension  $N_n$ :

**Assumption A:** Given  $\gamma \in \{-1/2, 0, 1/2\}$ ,  $a > 1$ ,  $C_1 > 0$  and differentiable functions  $f_1(x), f_2(x)$  with bounded derivatives  $f'_1(x), f'_2(x)$  such that

$$\sup_x \left| \mathbb{P}(m^\gamma T_m \leq x) - \Phi(x) - \frac{f_1(x)}{\sqrt{m}} - \frac{f_2(x)}{m} \right| \leq \frac{C_1}{m^a} \quad \text{for all } m \in \mathbb{N}. \quad (17)$$

*Remark 1.* Statistics satisfying Assumption A are shown in (14), (15) and (16).

**Assumption B:** Given constants  $b > 0$  and  $C_2 > 0$ , real numbers  $g_n$  with  $0 < g_n \uparrow \infty$  if  $n \rightarrow \infty$ , a distribution function  $H(y)$  with  $H(0+) = 0$  and a function  $h_2(y)$  of bounded variation that

$$\sup_{y \geq 0} \left| \mathbb{P} \left( \frac{N_n}{g_n} \leq y \right) - H(y) - \frac{h_2(y) \mathbb{I}_{\{b > 1\}}(b)}{n} \right| \leq \frac{C_2}{n^b} \quad \text{for all } n \geq 1. \quad (18)$$

where  $\mathbb{I}_A(x) = \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases}$  defines the indicator function of a set  $A \subset \mathbb{R}$ .

*Remark 2.* The random dimensions  $N_n(r)$  and  $N_n(s)$  given in (1) and (2), respectively, fulfill Assumption B as shown in [9, Propositions 1 and 2], see (29) and (39) below.

**Proposition 2.** Let  $\gamma \in \{1/2, 0, -1/2\}$  and both Assumption A and B as well as the following requirements on  $H(\cdot)$  and  $h_2(\cdot)$  are fulfilled

$$\left. \begin{aligned} i : & \quad H(1/g_n) \leq c_1 g_n^{-b} \quad \text{for } b > 0, \\ ii : & \quad \int_0^{1/g_n} y^{-1/2} dH(y) \leq c_2 g_n^{-b+1/2} \quad \text{for } b > 1/2, \\ iii : & \quad \int_0^{1/g_n} y^{-1} dH(y) \leq c_3 g_n^{-b+1} \quad \text{for } b > 1, \end{aligned} \right\} \quad (19)$$

$$\left. \begin{aligned} i : & \quad h_2(0) = 0, \quad \text{and } |h_2(1/g_n)| \leq c_4 n g_n^{-b} \quad \text{for } b > 1, \\ ii : & \quad \int_0^{1/g_n} y^{-1} |h_2(y)| dy \leq c_5 n g_n^{-b} \quad \text{for } b > 1, \end{aligned} \right\} \quad (20)$$

where  $b$  is the convergence rate in (18). Then for all  $n \geq 1$  is valid:

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P} \left( g_n^\gamma T_{N_n} \leq x \right) - G_{n,2}(x) \right| \leq C_1 \mathbb{E} (N_n^{-a}) + (C_3 D_n + C_4) n^{-b} + I_n, \quad (21)$$

where with  $a > 1, b > 0, f_1(z), f_2(z), h_2(y)$  are given in (17) and (18)

$$G_{n,2}(x) = \left\{ \begin{array}{ll} \int_0^\infty \Phi(xy^\gamma) dH(y), & 0 < b \leq 1/2, \\ \int_0^\infty \left( \Phi(xy^\gamma) + \frac{f_1(xy^\gamma)}{\sqrt{g_n y}} \right) dH(y) =: G_{n,1}(x), & 1/2 < b \leq 1, \\ G_{n,1}(x) + \int_0^\infty \frac{f_2(xy^\gamma)}{g_n y} dH(y) + \int_0^\infty \frac{\Phi(xy^\gamma)}{n} dh_2(y), & b > 1, \end{array} \right. \quad (22)$$

$$D_n = \sup_x \int_{1/g_n}^\infty \left| \frac{\partial}{\partial y} \left( \Phi(xy^\gamma) + \frac{f_1(xy^\gamma)}{\sqrt{g_n y}} + \frac{f_2(xy^\gamma)}{y g_n} \right) \right| dy, \quad (23)$$

$$I_n = \sup_x (|I_1(x, n)| + |I_2(x, n)|), \quad (24)$$

$$I_1(x, n) = \int_{1/g_n}^\infty \left( \frac{f_1(xy^\gamma) \mathbb{I}_{(0,1/2]}(b)}{\sqrt{g_n y}} + \frac{f_2(xy^\gamma)}{g_n y} \right) dH(y), \quad b \leq 1, \quad (25)$$

and

$$I_2(x, n) = \int_{1/g_n}^\infty \left( \frac{f_1(xy^\gamma)}{n \sqrt{g_n y}} + \frac{f_2(xy^\gamma)}{n g_n y} \right) dh_2(y), \quad b > 1. \quad (26)$$

The constants  $C_1, C_3, C_4$  are independent of  $n$ .

*Proof.* The proof is based on the statement in [2, Theorem 3.1] for  $\gamma \geq 0$ . Since in Theorems 1 and 2 in the present paper as well as in Christoph and Ulyanov [9, Theorems 1 and 2] the case  $\gamma = -1/2$  is also considered, therefore the proof was adapted to  $\gamma \in \{1/2, 0, -1/2\}$  in [9]. The conditions (19) and (20) guarantee integration range  $(0, \infty)$  of the integrals in (22). The approximation function  $G_{n,2}(x)$  in (22) is now a polynomial in  $g_n^{-1/2}$  and  $n^{-1/2}$ . Present Proposition 2 differs from Theorems 1 and 2 in [9] only by the term  $f_1(xy^\gamma) (g_n y)^{-1/2}$  and the added condition (19ii) to estimate this term. Therefore here the details are omitted.  $\square$

*Remark 3.* The domain  $[1/g_n, \infty)$  of integration depends on  $g_n$  in (23), (25) and (26). Some of the integrals in (25) and (26) could tend to infinity with  $1/g_n \rightarrow 0$  as  $n \rightarrow \infty$  and thus worsen the convergence rates of the corresponding terms. See (47) in Sect. 6.

In the next two sections we consider the statistics  $Z_m$  and  $Z_m^*$  defined in (10) and the cases when the random dimension  $N_n$  is given in either (1) or (2). We use Proposition 2 when the limit distributions of scaled statistics  $Z_{N_n}$  are scale mixtures  $G_\gamma(x) = \int_0^\infty \Phi(xy^\gamma) dH(y)$  with  $\gamma \in \{1/2, 0, -1/2\}$  that can be expressed in terms of the well-known distributions. We obtain non-asymptotic results for the statistics  $Z_{N_n}$  and  $Z_{N_n}^*$ , using second order approximations the statistics  $Z_m$  and  $Z_m^*$  given in (14) as well as for the random sample size  $N_n$ . In both cases the jumps of the distribution function of the random sample size  $N_n$  only affect the function  $h_2(y)$  in formula (18).

### 4 The Random Dimension $N_n(r)$ is Negative Binomial Distributed

The negative binomial distributed dimension  $N_n(r)$  has probability mass function (1)) and  $g_n = \mathbb{E}(N_n(r)) = r(n - 1) + 1$ . Schluter and Trede [21] (Sect. 2.1) underline the advantage of this distribution compared to the Poisson distribution for counting processes. They showed in a general unifying framework

$$\lim_{n \rightarrow \infty} \sup_y |\mathbb{P}(N_n(r)/g_n \leq y) - G_{r,r}(y)| = 0, \tag{27}$$

where  $G_{r,r}(y)$  is the Gamma distribution function with the identical shape and scale parameters  $r > 0$  and density

$$g_{r,r}(y) = \frac{r^r}{\Gamma(r)} y^{r-1} e^{-ry} \mathbb{I}_{(0, \infty)}(y) \quad \text{for all } y \in \mathbb{R}. \tag{28}$$

Statement (27) was proved earlier in Bening and Korolev [3, Lemma 2.2].

In [9, Proposition 1] the following inequality was proved for  $r > 0$ :

$$\sup_{y \geq 0} \left| \mathbb{P} \left( \frac{N_n(r)}{g_n} \leq y \right) - G_{r,r}(y) - \frac{h_{2;r}(y) \mathbb{I}_{\{r > 1\}}(r)}{n} \right| \leq \frac{C_2(r)}{n^{\min\{r, 2\}}}, \tag{29}$$

where  $h_{2;r}(y) = \frac{1}{2^r} g_{r,r}(y) ((y - 1)(2 - r) + 2Q_1(g_n y))$  for  $r > 1$ ,

$$Q_1(y) = 1/2 - (y - [y]) \quad \text{and } [y] \text{ is the integer part of a value } y. \tag{30}$$

Both Bening, Galieva and Korolev [2] and Gavrilenko, Zubov and Korolev [13] showed the rate of convergence in (29) for  $r \leq 1$ . In Christoph, Monakhov and Ulyanov [7, Theorem 1] the Chebyshev-Edgeworth expansion (29) for  $r > 1$  is proved.

*Remark 4.* The random dimension  $N_n(r)$  satisfies Assumption 2 of the Transfer Propositions 2 with  $g_n = \mathbb{E}N_n(r)$ ,  $H(y) = G_{r,r}(y)$ ,  $h_2(y) = h_{2;r}(y)$  and  $b = 2$ .

In (21), negative moment  $\mathbb{E}(N_n(r))^{-a}$  is required where  $m^{-a}$  is rate of convergence of Chebyshev-Edgeworth expansion for  $T_m$  in (17). Negative moments  $\mathbb{E}(N_n(r))^{-a}$  fulfill the estimate:

$$\mathbb{E}(N_n(r))^{-a} \leq C(a, r) \begin{cases} n^{-\min\{r, a\}}, r \neq a \\ \ln(n) n^{-a}, r = a \end{cases} \quad \text{for all } r > 0 \quad \text{and } a > 0. \tag{31}$$

For  $r = a$  the factor  $\ln n$  cannot be removed. In Christoph, Ulyanov and Bening [10, Corollary 4.2] leading terms for the negative moments of  $\mathbb{E}(N_n(r))^{-p}$  were derived for any  $p > 0$  that lead to (31).

The expansions of the length of the vector  $Z_m$  in (14) as well as of the sample correlation coefficient  $R_n$  in (15) and the angle  $\theta_m$  in (16) have as limit



distribution the standard normal distribution  $\Phi(x)$ . Therefore, with  $g_n = \mathbb{E}N_n(r)$  and  $\gamma \in \{1/2, 0, -1/2\}$ , limit distributions for

$$\mathbb{P}\left(g_n^\gamma(N_n(r))^{1/2-\gamma} Z_{N_n(r)} \leq x\right) \quad \text{are} \quad G_\gamma(x, r) = \int_0^\infty \Phi(xy^\gamma) dG_{r,r}(y).$$

These scale mixtures distributions  $G_\gamma(x, r)$  are calculated in Christoph and Ulyanov [9, Theorems 3–5]. We apply Proposition 2 to the statistics

$$T_{N_n(r)} = N_n(r)^{1/2-\gamma} Z_{N_n(r)} \quad \text{with the normalizing factor} \quad g_n^\gamma = \mathbb{E}(N_n(r))^\gamma.$$

The limit distributions are:

- for  $\gamma = 1/2$  and  $r > 0$  the **Student’s t-distribution**  $S_{2r}(x)$  with density

$$s_{2r}(x) = \frac{\Gamma(r + 1/2)}{\sqrt{2r\pi} \Gamma(r)} \left(1 + \frac{x^2}{2r}\right)^{-(r+1/2)}, \quad x \in \mathbb{R}, \quad (32)$$

- for  $\gamma = 0$  the **normal law**  $\Phi(x)$ ,
- for  $\gamma = -1/2$  and  $r = 2$  the **generalized Laplace distributions**  $L_2(x)$  with density  $l_2(x)$ :

$$L_2(x) = \frac{1}{2} + \frac{1}{2} \text{sign}(x) (1 - (1 + |x|) e^{-2|x|}) \quad \text{and} \quad l_2(x) = \left(\frac{1}{2} + |x|\right) e^{-2|x|}.$$

For arbitrary  $r > 0$  Macdonald functions  $K_{r-1/2}(x)$  occur in the density  $l_r(x)$ , which can be calculated in closed form for integer values of  $r$ .

The standard Laplace density with variance 1 is  $l_1(x) = \frac{1}{\sqrt{2}} e^{-\sqrt{2}|x|}$ .

**Theorem 1.** *Let  $Z_m$  and  $N_n(r)$  with  $r > 0$  be defined by (10) and (1), respectively. Suppose that (14) is satisfied for  $Z_m$  and (29) for  $N_n(r)$ . Then the following statements hold for all  $n \in \mathbb{N}_+$ :*

(i) **Student’s t approximation** using scaling factor  $\sqrt{\mathbb{E}N_n(r)}$  by  $Z_{N_n(r)}$

$$\sup_x \left| \mathbb{P}\left(\sqrt{g_n} Z_{N_n(r)} \leq x\right) - S_{2r;n}(x) \right| \leq C_r \begin{cases} n^{-\min\{r, 3/2\}}, & r \neq 3/2, \\ \ln(n) n^{-3/2}, & r = 3/2, \end{cases} \quad (33)$$

where

$$S_{2r;n}(x) = S_{2r}(x) + s_{2r}(x) \left( \frac{\sqrt{2}((2r-5)x^2 - 8r)}{12(2r-1)\sqrt{g_n}} \mathbb{I}_{\{r>1/2\}}(r) + \frac{96r^2x + (-64r^2 + 128r)x^3 + (4r^2 - 32r + 39)x^5}{(x^2 + 2r)(2r-1)g_n} \mathbb{I}_{\{r>1\}}(r) \right), \quad (34)$$

(ii) **Normal approximation** with random scaling factor  $N_n(r)$  by  $Z_{N_n(r)}$

$$\sup_x \left| \mathbb{P}(\sqrt{N_n(r)} Z_{N_n(r)} \leq x) - \Phi_{n,2}(x) \right| \leq C_r \begin{cases} n^{-\min\{r,3/2\}}, & r \neq 3/2, \\ \ln(n) n^{-3/2}, & r = 3/2, \end{cases} \quad (35)$$

where

$$\begin{aligned} \Phi_{n,2}(x) &= \Phi(x) + \frac{\sqrt{2} r \Gamma(r - 1/2)}{12 \Gamma(r) \sqrt{g_n}} (x^2 - 4) \varphi(x) \mathbb{I}_{\{r>1/2\}}(r) \\ &\quad + \frac{x^5 - 16x^3 + 24x}{144 g_n} \left( \frac{r}{r-1} \mathbb{I}_{\{r>1\}}(r) + \ln n \mathbb{I}_{\{r=1\}}(r) \right). \end{aligned} \quad (36)$$

(iii) **Generalized Laplace approximation** if  $r = 2$  with mixed scaling factor  $g_n^{-1/2} N_n(2)$  by  $Z_{N_n(2)}$

$$\sup_x \left| \mathbb{P} \left( g_n^{-1/2} N_n(2) Z_{N_n(2)} \leq x \right) - L_{n;2}(x) \right| \leq C_2 n^{-3/2} \quad (37)$$

where

$$\begin{aligned} L_{n;2}(x) &= L_2(x) - \frac{1}{3\sqrt{g_n}} \left( \frac{1}{\sqrt{2}} + \sqrt{2}|x| - x^2 \right) e^{-2|x|} \\ &\quad + \frac{1}{33 g_n} (12\sqrt{2}x - 15|x|x + 2x^3) e^{-2|x|}. \end{aligned} \quad (38)$$

## 5 The Random Dimension $N_n(s)$ is Discrete Pareto-Like Distributed

The Pareto-like distributed dimension  $N_n(s)$  has probability mass function (2) and  $\mathbb{E}(N_n(s)) = \infty$ . Hence  $g_n = n$  is chosen as normalizing sequence for  $N_n(s)$ .

Bening and Korolev [4, Sect. 4.3] showed that for integer  $s \geq 1$

$$\lim_{n \rightarrow \infty} \sup_{y>0} |\mathbb{P}(N_n(s) \leq ny) - H_s(y)| = 0.$$

where  $H_s(y) = e^{-s/y} \mathbb{I}_{(0, \infty)}(y)$  is the continuous distribution function of the inverse exponential  $W(s) = 1/V(s)$  with exponentially distributed  $V(s)$  having rate parameter  $s > 0$ . As  $\mathbb{P}(N_n(s) \leq y)$ , so  $H_s(y)$  is heavy tailed with shape parameter 1 and  $\mathbb{E}W(s) = \infty$ .

Lyamin [18] proved a bound  $|\mathbb{P}(N_n(s) \leq ny) - H_s(y)| \leq C/n$  and  $C < 0.37$  for integer  $s \geq 1$ .

In [9, Proposition 2] the following results are presented for  $s > 0$ :

$$\sup_{y>0} \left| \mathbb{P} \left( \frac{N_n(s)}{n} \leq y \right) - H_s(y) - \frac{h_{2;s}(y)}{n} \right| \leq \frac{C_3(s)}{n^2}, \quad \text{for all } n \in \mathbb{N}_+, \quad (39)$$

with  $H_s(y) = e^{-s/y}$  and  $h_{2;s}(y) = s e^{-s/y} (s - 1 + 2Q_1(ny)) / (2y^2)$  for  $y > 0$ , where  $Q_1(y)$  is defined in (30). Moreover

$$\mathbb{E}(N_n(s))^{-p} \leq C(p) n^{-\min\{p,2\}}, \quad (40)$$

where for  $0 < p \leq 2$  the order of the bound is optimal.

The Chebyshev-Edgeworth expansion (39) is proved in Christoph, Monakhov and Ulyanov [7, Theorem 4]. The leading terms for the negative moments  $\mathbb{E}(N_n(s))^{-p}$  were derived in Christoph, Ulyanov and Bening [10, Corollary 5.2] that lead to (40).

*Remark 5.* The random dimension  $N_n(s)$  satisfies Assumption 2 of the Transfer Propositions 2 with  $H_s(y) = e^{-s/y}$ ,  $h_2(y) = h_{2;s}(y)$ ,  $g_n = n$  and  $b = 2$ .

With  $g_n = n$  and  $\gamma \in \{1/2, 0, -1/2\}$ , the limit distributions for

$$\mathbb{P}\left(n^\gamma N_n(s)^{1/2-\gamma} Z_{N_n(s)} \leq x\right) \quad \text{are now} \quad G_\gamma(x, s) = \int_0^\infty \Phi(xy^\gamma) dH_s(y).$$

These scale mixtures distributions  $G_\gamma(x, s)$  are calculated in Christoph and Ulyanov [9, Theorems 6–8]. We apply Proposition 2 to statistics

$$T_{N_n(s)} = N_n(s)^{1/2-\gamma} Z_{N_n(s)} \quad \text{with the normalizing factor } n^\gamma.$$

The limit distributions are:

- for  $\gamma = 1/2$  **Laplace distributions**  $L_{1/\sqrt{s}}(x)$  with density

$$l_{1/\sqrt{s}}(x) = \sqrt{s/2} e^{-\sqrt{2s}|x|},$$

- for  $\gamma = 0$  the **standard normal law**  $\Phi(x)$  and
- for  $\gamma = -1/2$  the **scaled Student’s t-distribution**  $S_2^*(x; \sqrt{s})$  with density

$$s_2^*(x; \sqrt{s}) = \frac{1}{2\sqrt{2s}} \left(1 + \frac{x^2}{2s}\right)^{-3/2}.$$

**Theorem 2.** *Let  $Z_m$  and  $N_n(s)$  with  $s > 0$  be defined by (10) and (2), respectively. Suppose that (14) is satisfied for  $Z_m$  and (39) for  $N_n(s)$ . Then the following statements hold for all  $n \in \mathbb{N}_+$ :*

- (i) **Laplace approximation** with non-random scaling factor  $n^\gamma$  by  $Z_{N_n(s)}$ :

$$\sup_x \left| \mathbb{P}\left(\sqrt{n} Z_{N_n(s)} \leq x\right) - L_{1/\sqrt{s};n}(x) \right| \leq C_\infty n^{-3/2} \tag{41}$$

where

$$\begin{aligned} L_{1/\sqrt{s};n}(x) = & L_{1/\sqrt{s}}(x) + l_{1/\sqrt{s}}(x) \left( \frac{\sqrt{2}}{12s\sqrt{n}} (sx^2 - 2(1 + \sqrt{2s}|x|)) \right. \\ & \left. + \frac{s}{72n} \left( \frac{x^3|x|}{\sqrt{2s}} - \frac{8x^2}{s} + \frac{6x}{s^2} (1 + \sqrt{2s}|x|) \right) \right) \end{aligned} \tag{42}$$

,

(ii) **Normal approximation** with random scaling factor  $\sqrt{N_n(s)}$  by  $Z_{N_n(r)}$ :

$$\sup_x \left| \mathbb{P} \left( \sqrt{N_n(s)} Z_{N_n(s)} \leq x \right) - \Phi_{n,2}(x) \right| \leq C_s n^{-3/2}, \tag{43}$$

where

$$\Phi_{n,2}(x) = \Phi(x) + \varphi(x) \left( \frac{\sqrt{2\pi}(x^2 - 4)}{24\sqrt{n}} + \frac{x^5 - 16x^3 + 24x}{144sn} \right) \tag{44}$$

(iii) **Scaled Student's t-distribution** with mixed scaling factor by  $Z_{N_n(s)}$

$$\sup_x \left| \mathbb{P} \left( n^{-1/2} N_n(s) Z_{N_n(s)} \leq x \right) - S_{n,2}^*(x) \right| \leq C_s n^{-3/2}, \tag{45}$$

where

$$S_{n,2}^*(x; \sqrt{s}) = S_2^*(x; \sqrt{s}) + s_2^*(x; \sqrt{s}) \left( -\frac{\sqrt{2}(x^2 + 8s)}{12(2s + x^2)\sqrt{n}} + \frac{1}{144n} \left( \frac{105x^5}{(2s + x^2)^3} + \frac{240x^3}{(2s + x^2)^2} + \frac{72x}{2s + x^2} \right) \right). \tag{46}$$

## 6 Proofs of Main Results

*Proof.* The proofs of Theorems 1 and 2 are based on Proposition 2. The structure of the functions  $f_1, f_2$  and  $h_2$  in Assumptions A and B is similar to the structure of the corresponding functions in Conditions 1 and 2 in [9]. Therefore, the estimates of the term  $D_n$  and of the integrals  $I_1(x, n)$  and  $I_2(x, n)$  in (23), (25) and (24) as well as the validity of (19) and (20) in Proposition 2 when  $H(y)$  is  $G_{r,r}(y)$  or  $H_s(y)$  can be shown analogously to the proofs for Lemmas 1, 2 or 4 in [9]. In Remark 3 above it was pointed out that the integrals in (25) and (25) can degrade the convergence rate. Let  $r < 1$ . With  $|f_2(xy^\gamma)| \leq c^*$  we get

$$\int_{1/g_n}^\infty \frac{|f_2(xy^\gamma)|}{g_n y} dG_{r,r}(y) \leq \frac{c^* r^r}{\Gamma(r) g_n} \int_{1/g_n}^\infty y^{r-2} dy \leq \frac{c^* r^r}{(1-r)\Gamma(r)} g_n^{-r}. \tag{47}$$

The additional term  $f_1(xy^\gamma)(g_n y)^{-1/2}$  in (17) in Assumption A is to be estimated with condition (19ii).

Moreover, the bounds for  $\mathbb{E}(N_n)^{-3/2}$  follow from (31) and (40), since  $a = 3/2$  in Assumption A, considering the approximation (14).

The integrals in (22) in Proposition 2 are still to be calculated. Similar integrals are calculated in great detail in the proofs of Theorems 3–8 in [9]. To obtain (34), we compute the integrals with Formula 2.3.3.1 in Prudnikov et al. [20]

$$M_\alpha(x) = \frac{r^r}{\Gamma(r)\sqrt{2\pi}} \int_0^\infty y^{\alpha-1} e^{-(r+x^2/2)y} dy = \frac{\Gamma(\alpha) r^{r-\alpha}}{\Gamma(r)\sqrt{2\pi}} \left( 1 + x^2/(2r) \right)^{-\alpha} \tag{48}$$

for  $\alpha = r - 1/2, r + 1/2, r + 3/2$  and  $p = r + x^2/2$ .

Lemma 2 in [9] and  $\int_0^\infty y^{-1} dG_{r,r}(y) = r/(r-1)$  for  $r > 1$  lead to (36).

To show (38) we use Formula 2.3.16.2 in [20] with  $n = 0, 1$  and Formula 2.3.16.3 in [20] with  $n = 1, 2$  and  $p = 2$  and  $q = x^2/2$ .

To obtain (42), we calculate the integrals again with Formula 2.3.16.3 in [20], with  $p = x^2/2 > 0$ ,  $q = s > 0$ ,  $n = 0, 1, 2$ .

Lemma 4 in [9] and  $\int_0^\infty y^{-a-1} e^{-s/y} dy = s^{-a} \Gamma(a)$  for  $a = 3/2, 2$  lead to (44).

Finally, in  $\int_0^\infty f_k(x/y^\gamma) y^{-2-k/2} e^{-s/y} dy$  we use the substitution  $s/y = u$  to obtain, with (48), the terms in (46).  $\square$

**Acknowledgements.** Theorem 1 has been obtained under support of the Ministry of Education and Science of the Russian Federation as part of the program of the Moscow Center for Fundamental and Applied Mathematics under the agreement № 075-15-2019-1621. Theorem 2 was proved within the framework of the HSE University Basic Research Program.

## References

1. Aoshima, M., Shen, D., Shen, H., Yata, K., Zhou, Y.-H., Marron, J.S.: A survey of high dimension low sample size asymptotics. *Aust. N. Z. J. Stat.* **60**(1), 4–19 (2018). <https://doi.org/10.1111/anzs.12212>
2. Bening, V.E., Galieva, N.K., Korolev, V.Y.: Asymptotic expansions for the distribution functions of statistics constructed from samples with random sizes [in Russian]. *Inf. Appl. IPI RAN* **7**(2), 75–83 (2013)
3. Bening, V.E., Korolev, V.Y.: On the use of Student’s distribution in problems of probability theory and mathematical statistics. *Theory Probab. Appl.* **49**(3), 377–391 (2005)
4. Bening, V.E., Korolev, V.Y.: Some statistical problems related to the Laplace distribution [in Russian]. *Inf. Appl. IPI RAN* **2**(2), 19–34 (2008)
5. Bobkov, S.G., Naumov, A.A., Ulyanov V.V.: Two-sided inequalities for the density function’s maximum of weighted sum of chi-square variables. [arXiv:2012.10747v1](https://arxiv.org/pdf/2012.10747v1) (2020). <https://arxiv.org/pdf/2012.10747.pdf>
6. Buddana, A., Kozubowski, T.J.: Discrete Pareto distributions. *Econ. Qual. Control* **29**(2), 143–156 (2014)
7. Christoph, G., Monakhov, M.M., Ulyanov, V.V.: Second-order Chebyshev-Edgeworth and Cornish-Fisher expansions for distributions of statistics constructed with respect to samples of random size. *J. Math. Sci. (N.Y.)* **244**(5), 811–839 (2020). Translated from *Zapiski Nauchnykh Seminarov POMI*, 466, *Veroyatnost i Statistika*. 26, 167–207 (2017)
8. Christoph, G., Prokhorov, Yu., Ulyanov, V.: On distribution of quadratic forms in Gaussian random variables. *Theory Prob. Appl.* **40**(2), 250–260 (1996)
9. Christoph, G., Ulyanov, V.V.: Second order expansions for high-dimension low-sample-size data statistics in random setting. *Mathematics* **8**(7), 1151 (2020)
10. Christoph, G., Ulyanov, V.V., Bening, V.E.: Second order expansions for sample median with random sample size. [arXiv:1905.07765v2](https://arxiv.org/pdf/1905.07765v2) (2020). <https://arxiv.org/pdf/1905.07765.pdf>

11. Christoph, G., Ulyanov, V.V., Fujikoshi, Y.: Accurate approximation of correlation coefficients by short Edgeworth-Chebyshev expansion and its statistical applications. In: Shiryaev, A.N., Varadhan, S.R.S., Presman, E.L. (eds.) *Prokhorov and Contemporary Probability Theory. In Honor of Yuri V. Prokhorov*. Springer Proceedings in Mathematics & Statistics, vol. 33, pp. 239–260. Springer, Heidelberg (2013). [https://doi.org/10.1007/978-3-642-33549-5\\_13](https://doi.org/10.1007/978-3-642-33549-5_13)
12. Fujikoshi, Y., Ulyanov, V.V., Shimizu, R.: *Multivariate Statistics. High-Dimensional and Large-Sample Approximations*. Wiley Series in Probability and Statistics. Wiley, Hoboken (2010)
13. Gavrilenko, S.V., Zubov, V.N., Korolev, V.Y.: The rate of convergence of the distributions of regular statistics constructed from samples with negatively binomially distributed random sizes to the Student distribution. *J. Math. Sci. (N.Y.)* **220**(6), 701–713 (2017)
14. Hall, P., Marron, J.S., Neeman, A.: Geometric representation of high dimension, low sample size data. *J. R. Stat. Soc. Ser.* **67**, 427–444 (2005)
15. Johnson, N.L., Kotz, S., Balakrishnan, N.: *Continuous Univariate Distributions*, vol. 2, 2nd edn. Wiley, New York (1995)
16. Kawaguchi, Y., Ulyanov, V.V., Fujikoshi, Y.: Asymptotic distributions of basic statistics in geometric representation for high-dimensional data and their error bounds (Russian). *Inf. Appl.* **4**, 12–17 (2010)
17. Konishi, S.: Asymptotic expansions for the distributions of functions of a correlation matrix. *J. Multivar. Anal.* **9**, 259–266 (1979)
18. Lyamin, O.O.: On the rate of convergence of the distributions of certain statistics to the Laplace distribution. *Mosc. Univ. Comput. Math. Cybern.* **34**(3), 126–134 (2010)
19. Petrov, V.V.: *Limit Theorems of Probability Theory. Sequences of Independent Random Variables*. Clarendon Press, Oxford (1995)
20. Prudnikov, A.P., Brychkov, Y.A., Marichev, O.I.: *Integrals and Series, Volume 1: Elementary Functions*, 3rd edn. Gordon & Breach Science Publishers, New York (1992)
21. Schluter, C., Trede, M.: Weak convergence to the student and Laplace distributions. *J. Appl. Probab.* **53**(1), 121–129 (2016)