# An Edgeworth expansion for finite-population *U*-statistics

MINDAUGAS BLOZNELIS[1] and FRIEDRICH GÖTZE[2]

[1]*Department of Mathematics, Vilnius University, Naugarduko 24, Vilnius 2006, Lithuania.*
*E-mail: mindaugas.bloznelis@maf.vu.lt*
[2]*Fakultät für Mathematik, Universität Bielefeld, Postfach 100131, 33501 Bielefeld, Germany.*
*E-mail: goetze@mathematik.uni-bielefeld.de*

Suppose that $U$ is a $U$-statistic of degree 2 based on $N$ random observations drawn without replacement from a finite population. For the distribution of a standardized version of $U$ we construct an Edgeworth expansion with remainder $O(N^{-1})$ provided that the linear part of the statistic satisfies a Cramér type condition.

*Keywords:* central limit theorem; Edgeworth expansion; finite population; $U$-statistic; sampling without replacement

## 1. Introduction and results

Let $\mathscr{A} = \{a_1, \ldots, a_n\}$ denote a population of size $n$ and let $\mathscr{H} : \mathscr{A} \times \mathscr{A} \to \mathbb{R}$ denote symmetric function of its two arguments. By $X_1, \ldots, X_N$, $N \leqslant n$, we denote random variables with values in $\mathscr{A}$ such that $X = \{X_1, \ldots, X_N\}$ represents a random sample from $\mathscr{A}$ of size $N$ drawn without replacement, i.e. $\mathrm{P}\{X = B\} = \binom{n}{N}^{-1}$ for any subset $B \subset \mathscr{A}$ of size $N$. We shall investigate the second-order asymptotics of the distribution of the statistic

$$U = \sum_{1 \leqslant i \leqslant j \leqslant N} \mathscr{H}(X_i, X_j).$$

We assume that the statistic is centred. Write

$$U = L + Q, \tag{1.1}$$

where

$$L = \sum_{i=1}^{N} g_1(X_i) \quad \text{and} \quad Q = \sum_{1 \leqslant i \leqslant j \leqslant N} g_2(X_i, X_j)$$

are respectively the linear and the quadratic part of the statistic. Here

$$g_1(x) = (N-1)t(x), \qquad t(x) = \frac{n-1}{n-2}\mathrm{E}(\mathscr{H}(X_1, X_2)|X_1 = x)$$

and

$$g_2(x_1, x_2) = \mathscr{H}(x_1, x_2) - t(x_1) - t(x_2).$$

Since

$$\mathrm{E}(g_2(X_1, X_2)|X_1 = x) = 0, \qquad \text{for all } x \in \mathscr{A}, \tag{1.2}$$

the random variables $g_1(X_i)$ and $g_2(X_j, X_k)$, $1 \leqslant i \leqslant N$, $1 \leqslant j < k \leqslant N$ (and thus $L$ and $Q$) are uncorrelated. If the linear part $L$ dominates the statistic, for large $N$, the distribution of $U$ can be approximated by a Gaussian distribution using the central limit theorem (CLT).

The asymptotic normality of linear statistics based on samples drawn without replacement from finite populations has been studied by a number of authors. Erdős and Rényi (1959) proved the CLT under very mild conditions. The rate of convergence in the CLT was first studied by Bikelis (1969). Berry–Esseen bounds of order $O(N^{-1/2})$ were obtained by Höglund (1978). Robinson (1978) proved the validity of an Edgeworth expansion with a remainder of order $O(N^{-3/2})$; see also Bickel and van Zwet (1978).

Nandi and Sen (1963) studied the asymptotic behaviour of finite-population $U$-statistics and showed that under proper regularity conditions the sequence of distributions of normalized $U$-statistics converges to the standard normal distribution. The rate of this convergence was investigated by Zhao and Chen (1987; 1990), Kokic and Weber (1990; 1991) and, as a particular case of the rate of convergence of general multivariate sampling statistics, by Bolthausen and Götze (1993). In the case of independent and identically distributed (i.i.d.) observations the second-order asymptotic theory has been developed for $U$-statistics: see Bickel (1974), Götze (1979), Callaert *et al.* (1980), Bickel *et al.* (1986) and, for more general asymptotically normal symmetric statistics, Bentkus *et al.* (1997). In contrast to the independent case, there are only a few results concerned with higher-order asymptotics of nonlinear finite population statistics. Babu and Singh (1985) proved the validity of an Edgeworth expansion with a remainder $o(N^{-1/2})$ for finite-population multivariate sample mean and applied this result to establish expansions for statistics that can be represented as smooth functions of multivariate sample means, e.g. Student's $t$. Kokic and Weber (1990) established a one-term Edgeworth expansion with the remainder $o(N^{-1/2})$ for finite-population $U$-statistics of degree 2.

By way of comparison to the results described above, we shall provide an *explicit* remainder term of order $O(N^{-1})$ for finite population $U$-statistics which is optimal assuming a Cramér condition on the linear term only. The proof is based on a finite-population variant of Hoeffding's decomposition as well as the Erdős–Rényi representation and some ideas due to Bentkus *et al.* (1997) such as data-dependent smoothing.

Assume that

$$\sigma^2 = N \mathrm{E} g_1^2(X_1) > 0.$$

The distribution function of the standardized statistic, $F(x) = \mathrm{P}\{U \leqslant x\sigma\}$, will be approximated by the one-term Edgeworth expansion,

$$G(x) = \Phi\left(\frac{x}{\sqrt{q}}\right) - \frac{(q - p)q^{-1/2}\alpha + 3q^{1/2}\kappa}{6\sigma^3 N^{1/2}}\Phi'''\left(\frac{x}{\sqrt{q}}\right). \tag{1.3}$$

Here $\Phi(x)$ is the standard normal distribution function,

$$p = N/n, \qquad q = 1 - p$$

and

$$\alpha = N^{3/2}\mathrm{E}g_1^3(X_1), \qquad \kappa = N^{5/2}\mathrm{E}g_2(X_1, X_2)g_1(X_1)g_1(X_2). \tag{1.4}$$

We shall derive bounds for the remainder

$$\Delta = \sup_{x \in \mathbb{R}}|F(x) - G(x)|.$$

To prove the validity of an Edgeworth expansion, i.e. to establish bounds for $\Delta$, in addition to moment conditions one needs to impose a smoothness condition, cf. Bickel and Robinson (1982). For instance, in the classical case of standardized sums $S = (Y_1 + \ldots + Y_N)/\sqrt{N}$ of i.i.d. random variables $Y_1, \ldots, Y_N$ such that $\mathrm{E}Y_1 = 0$, $\mathrm{E}Y_1^2 = 1$ and $\mathrm{E}Y_1^4 < \infty$, asymptotic expansions for the distribution $F_S$ of $S$ with remainder $O(N^{-1})$ are obtained assuming Cramér's condition,

$$\sup_{|t| > a}|\mathrm{E}\exp\{itY_1\}| < 1. \tag{C}$$

Bentkus *et al.* (1997) introduced a local version of Cramér's condition (**C**), namely,

$$\rho_{Y_1}(a, b) := 1 - \sup_{a \leqslant |t| \leqslant b}|\mathrm{E}\exp\{itY_1\}| > 0. \tag{C'}$$

Condition (**C'**) (with $a = 1/\mathrm{E}|Y_1|^3$ and $b = N^{1/2}$) is somewhat weaker than (**C**) but still sufficient to prove the validity of Edgeworth expansions for $F_S$ up to order $O(N^{-1})$. This modification is useful in more general situations, where $Y_1$ depends on $N$ in an implicit way; see Bentkus *et al.* (1997).

For a sufficiently small absolute constant $b_1$, say $b_1 = 0.0001$, we shall assume that the distribution of the random variable $Z = \sqrt{N}g_1(X_1)/\sigma$ satisfies condition (**C'**) with $a' = b_1/\mathrm{E}|Z|^3$ and $b' = N^{1/2}$, i.e.

$$\rho = \rho_Z(a', b') > 0. \tag{1.5}$$

Write, for $r = 1, 2, \ldots,$

$$\beta_r = \mathrm{E}|N^{1/2}g_1(X_1)|^r \quad \text{and} \quad \gamma_r = \mathrm{E}|N^{3/2}g_2(X_1, X_2)|^r. \tag{1.6}$$

Then the following estimate holds for the remainder $\Delta$.

**Theorem 1.1.** *There exists an absolute constant $A > 0$ such that*

$$\sup_{x \in \mathbb{R}}|F(x) - G(x)| \leqslant \frac{A}{N}\frac{\beta_4 + \gamma_4}{\rho^2 q^2 \sigma^4}.$$

For linear statistics we obtain the following result.

**Theorem 1.2.** *There exists an absolute constant $B > 0$ such that*

$$\left| \mathrm{P}\{L \leqslant x\} - \Phi\left(\frac{x}{\sqrt{q}}\right) + \frac{(q-p)q^{-1/2}\alpha}{6\sigma_3 N^{1/2}} \Phi'''\left(\frac{x}{\sqrt{q}}\right) \right| \leqslant \frac{B}{N} \frac{\beta_4}{\rho^2 q \sigma^4}.$$

The estimates in Theorems 1.1 and 1.2 hold for any fixed sample size $N$, population size $n$ and functions $\mathcal{H}$. If $\beta_4/\sigma^4$ and $\gamma_4/\sigma^4$ are bounded and $q$ and $\rho$ are bounded away from 0 and $N \to \infty$ and $n \to \infty$, then these results establish Edgeworth expansions with the remainder $O(N^{-1})$.

The case where $n \to \infty$ and $N$ is fixed corresponds to the i.i.d. situation. By the law of large numbers we obtain a corollary for independent observations. Let $\mathcal{E}$ denote a measurable space and let $\mathcal{X}_1, \mathcal{X}_2, \ldots$ be i.i.d. random variables with values in $\mathcal{E}$. Write

$$\tilde{U} = \sum_{1 \leqslant i \leqslant j \leqslant N} \mathcal{H}(\mathcal{X}_i, \mathcal{X}_j).$$

Here $\mathcal{H} : \mathcal{E} \times \mathcal{E} \to \mathbb{R}$ denotes a measurable function symmetric in its two arguments such that $\mathrm{E}\mathcal{H}^2(\mathcal{X}_1, \mathcal{X}_2) < \infty$. We assume that $\mathrm{E}\tilde{U} = 0$ and decompose

$$\tilde{U} = \sum_{i=1}^{N} \tilde{g}_1(X_i) + \sum_{1 \leqslant i \leqslant j \leqslant N} \tilde{g}_2(X_i, X_j).$$

Here $\tilde{g}_1$ and $\tilde{g}_2$ are defined in the same way as $g_1$ and $g_2$, but using $\tilde{t}(x) = \mathrm{E}(\mathcal{H}(\mathcal{X}_1, \mathcal{X}_2)|\mathcal{X}_1 = x)$ instead of $t(x)$. Let $\tilde{\sigma}$, $\tilde{\alpha}$, $\tilde{\beta}_k$, $\tilde{\gamma}_k$, $k = 2, 3, 4$, and $\tilde{\kappa}$ denote the moments of $\tilde{g}_1(\mathcal{X}_1)$ and $\tilde{g}_2(\mathcal{X}_1, \mathcal{X}_2)$ corresponding to $\sigma$, $\alpha$, $\beta_k$, $\gamma_k$ and $\kappa$. We shall assume that

$$\tilde{\rho} = \rho_{\mathcal{Z}}(\tilde{a}, \tilde{b}) > 0, \qquad \mathcal{Z} = \sqrt{N}\tilde{g}_1(\mathcal{X}_1)/\tilde{\sigma},$$

where $\tilde{a} = b_1/\mathrm{E}|\mathcal{Z}|^3$ and $\tilde{b} = \sqrt{N}$. Then we have:

**Corollary 1.3.** *There exists an absolute constant $A > 0$ such that*

$$\left| \mathrm{P}\{\tilde{U} \leqslant \tilde{\sigma}x\} - \Phi(x) + \frac{\tilde{\alpha} + 3\tilde{\kappa}}{6\tilde{\sigma}^3 N^{1/2}} \Phi'''(x) \right| \leqslant \frac{A}{N} \frac{\tilde{\beta}_4 + \tilde{\gamma}_4}{\tilde{\rho}^2 \tilde{\sigma}^4}.$$

Hence Theorem 1.1, which yields this result as a special case, may be regarded as a partial extension of the result of Bentkus *et al.* (1997) to a simple random sampling model. They proved the validity of an Edgeworth expansion with remainder $O(N^{-1})$ for general symmetric asymptotically normal statistics based on i.i.d. observations. In the case of $U$-statistics of degree 2 their result yields the estimate as in Corollary 1.3 but with a lower moment $\tilde{\gamma}_3/\tilde{\sigma}^3$ instead of $\tilde{\gamma}_4/\tilde{\sigma}^4$ in the remainder.

An example given in Theorem 1.4 of Bentkus *et al.* (1997) shows that a Cramér type condition on the linear part and the existence of moments of arbitrarily high order of the linear and quadratic parts of the statistic (based on i.i.d. observations) are not sufficient to obtain higher-order approximations (those with remainders $o(N^{-1})$) to the distribution function of $U$. Hence, in this sense Corollary 1.3 and thus Theorem 1.1 are the best possible. To prove the validity of an Edgeworth expansion with remainder $o(N^{-1})$, one

needs in addition to impose a smoothness condition on the distribution of the quadratic part; see, for example, Bickel *et al.* (1986).

Let us compare our results with those of Robinson (1978) and Kokic and Weber (1990). Robinson (1978) proved the validity of a two-term Edgeworth expansion with remainder $O(N^{-3/2})$ for linear statistics like $L$ in (1.1) assuming the following Cramér type condition. This condition, first used in Albers *et al.* (1976), requires for a random variable $Z$ that there exists an $\varepsilon > 0$ such that

$$\tau_Z(\varepsilon, a, b) = 1 - \sup_{s \in \mathbb{R}, a \leq |t| \leq b} \mathrm{P}\{tZ \in \mathscr{L}^\varepsilon + s\} > 0. \tag{c}$$

Here $\mathscr{L} = \{2\pi r, r = 0, \pm 1, \pm 2, \ldots\}$ and $\mathscr{B}^\varepsilon$ denotes the $\varepsilon$-neighbourhood of a set $\mathscr{B} \subset \mathbb{R}$. Notice that $\varepsilon_1 \leq \varepsilon_2$ implies $\tau_Z(\varepsilon_1, a, b) \geq \tau_Z(\varepsilon_2, a, b)$. Robinson assumed that given $C' > 0$

$$\text{there exist } \varepsilon, \delta > 0 \text{ and } C > 0 \text{ such that } \tau_Z(\varepsilon, a, b) > \delta, \tag{1.7}$$

for

$$Z = \sqrt{N} g_1(X_1)/\sigma, \qquad a^{-1} = \max_{1 \leq i \leq n} |z_i|/C', \qquad b^{-1} = p\mathrm{E}|Z|^5/(CN).$$

Here $\{z_1, \ldots, z_n\}$ denotes the set of values of the random variable $Z$. Note that $\max_i|z_i| = \max_i|z_i|\mathrm{E}Z^2 \geq \mathrm{E}|Z|^3$, because of $\mathrm{E}Z^2 = 1$. For a sequence of finite-population linear statistics, say $(L_n)$, Robinson's (1978) theorem establishes an Edgeworth expansion with remainder $O(N^{-3/2})$ provided that $\beta_5/\sigma^5$ is bounded, $p$ and $q$ are bounded away from 0 and (1.7) holds with $\varepsilon$, $\delta$ and $C$ not depending on $n$ as $n \to \infty$. Robinson's (1978) result was used by Kokic and Weber (1990) to show $\Delta = o(N^{-1/2})$. The bounds for the remainders in these papers involve constants which implicitly depend on $p$.

In Section 2 we compare conditions (**c**) and (1.5). Proofs of Theorems 1.1 and 1.2 and of Corollary 1.3 are given in Sections 3 and 4. Auxiliary results are gathered together in Section 5.

## 2. Smoothness conditions

Modifications of Cramér's condition (**C**) that ensure the validity of Edgeworth expansions for sums of random variables assuming a finite number of values only were considered by Albers *et al.* (1976), van Zwet (1982), Does (1983) and Schneller (1989); see also Bickel and Robinson (1982). In this section we show that a Cramér type condition used in Albers *et al.* (1976) and Robinson (1978) is equivalent to that introduced in Bentkus *et al.* (1997) – namely, that the conditions (1.5) and (**c**) are equivalent. More specifically, given a random variable $Z$ and numbers $0 < a < b$, (1.5) implies $\tau_Z(\rho, a, b) > \rho/4$. Furthermore, if (**c**) holds for some $\varepsilon > 0$, then $\rho_Z(a, b) > \varepsilon^2 \tau_Z(\varepsilon, a, b)/\pi^2$; see Lemma 2.1 below.

In order to check condition (**c**) one needs to maximize a bivariate function over the set $(s, t) \in [-\pi, \pi] \times \{a \leq |t| \leq b\}$. Such a (maximization) problem can be difficult to solve numerically. A symmetrization argument suggests a version of condition (**c**) which is easier to check. Let $Z'$ denote an independent copy of $Z$ and let $Z^* = Z - Z'$ denote a symmetrization of $Z$. The condition

there exists $\varepsilon > 0$     such that $\tau_Z^*(\varepsilon, a, b) = 1 - \sup\limits_{a \leqslant |t| \leqslant b} \mathrm{P}\{tZ^* \in \mathscr{L}^\varepsilon\} > 0$      $(\mathbf{c}^*)$

requires the estimation of the maximum of an univariate function only. Condition $(\mathbf{c}^*)$ was proposed by V. Bentkus. Notice that $\varepsilon_1 < \varepsilon_2$ implies $\tau_Z^*(\varepsilon_2, a, b) \leqslant \tau_Z^*(\varepsilon_1, a, b)$. The following lemma shows that conditions $(\mathbf{c}^*)$ and $(\mathbf{c})$ are equivalent. Write

$$\delta_Z(a, b) = 1 - \sup\{\mathrm{E}\cos(tZ + s) : s \in \mathbb{R}, \ a \leqslant |t| \leqslant b\}. \tag{2.1}$$

**Lemma 2.1** *Let $Z$ be a random variable. For $0 < a < b$ and $0 < \varepsilon < \pi$, write*

$$\rho = \rho_Z(a, b), \qquad \tau_\varepsilon = \tau_Z(\varepsilon, a, b) \qquad \tau_\varepsilon^* = \tau_Z^*(\varepsilon, a, b), \qquad u = \pi^{-1}\varepsilon\tau_\varepsilon^*, \qquad v = \pi^{-1}\varepsilon\tau_\varepsilon.$$

*The following inequalities hold:*

$$\frac{\varepsilon^2 \tau_\varepsilon}{\pi^2} \leqslant \rho \leqslant 4\tau_\rho, \qquad \frac{\varepsilon^2 \tau_\varepsilon^*}{\pi^2} \leqslant \rho \leqslant 4\tau_\rho^*, \qquad \tau_v^* \geqslant \frac{\varepsilon^2 \tau_\varepsilon}{2\pi^2}, \qquad \tau_u \geqslant \frac{\varepsilon^2 \tau_\varepsilon^*}{4\pi^2}, \qquad \delta_Z(a, b) \geqslant \rho.$$

The proof of Lemma 2.1 is elementary; see Bloznelis and Götze (1997).

# 3. Proofs

Throughout this section and the next we shall assume without loss of generality that $\beta_2 = 1$. Since the proof of our main result, Theorem 1.1 is rather complex and involved we shall first outline the various steps.

In the first step, choosing $m \approx \ln N$, we replace the statistic $U$ by

$$U_1 = L' + U', \qquad U' = g_1(X_{m+1}) + \ldots + g_1(X_N) + \sum_{m+1 \leqslant i < j \leqslant N} g_2(X_i, X_j), \tag{3.1}$$

where

$$L' = l(X_1) + \ldots + l(X_m),$$

with

$$l(x) = g_1(x) + l_0(x), \qquad l_0(x) = \sum_{j=m+1}^{N} g_2(x, X_j),$$

is a conditionally linear statistic given $X_{m+1}, \ldots, X_N$. Write

$$F_X(x) = \mathrm{P}\{U_1 \leqslant x | X_{m+1}, \ldots, X_N\}, \qquad f_1(t) = \mathrm{E}(\exp\{itU_1\}|X_{m+1}, \ldots, X_N).$$

In the second step we construct upper/lower bounds for conditional probabilities

$$F_X(x+) \leqslant \frac{1}{2} + \mathrm{VP}\int_{\mathbb{R}} \exp\{-ixt\} \frac{1}{H} K\left(\frac{t}{H}\right) f_1(t)\, dt,$$

$$F_X(x-) \geqslant \frac{1}{2} - \mathrm{VP}\int_{\mathbb{R}} \exp\{-ixt\} \frac{1}{H} K\left(\frac{-t}{H}\right) f_1(t)\, dt,$$

where $F(x+) = \lim_{z \downarrow x} F(z)$, $F(x-) = \lim_{z \uparrow x} F(z)$ and VP denotes Cauchy's principal value (Prawitz's (1972) smoothing lemma). The bounded weight function $K(t/H)$, vanishing for $|t| > H$, and the cut-off $H = O(N)$ are specified below. Taking expectations of the left- and right-hand sides respectively, we obtain upper and lower bounds for the distribution function $F_1(x) = P\{U_1 \leqslant x\}$; see (3.7) and (3.8) below.

In the third step we construct a bound for the integral of $f_1(t)K(t/H)$ over the region $cN^{1/2} \leqslant |t| \leqslant H$. In the classical linear statistic case the bounds for the characteristic function for large values of $t$, like $cN^{1/2} \leqslant |t| \leqslant CN$, are implied by Cramér's condition (**C**). We write

$$|f_1(t)| \leqslant |E(\exp\{it(l(X_1) + \ldots + l(X_m))\}|X_{m+1}, \ldots, X_N)|$$

and show that the Cramér condition $|E \exp\{itg_1(X_1)\}| < 1 - \rho$ (we do not require $|E \exp\{itl(X_1)\}| < 1 - \rho$), in combination with a suitable choice of the cut-off $H = H(X_{m+1}, \ldots, X_N)$ implies a bound like $|f_1| \leqslant (1 - c\rho)^m$, for some $0 < c < 1$. The techniques are somewhat complicated by the fact that $X_1, \ldots, X_m$ are exchangeable only and we get the independence via the Erdős-Rényi decomposition for (conditional and unconditional) characteristic functions.

In the next step we interchange the conditional characteristic function with the unconditional one by changing the order of integration with respect to Lebesgue measure and with respect to the distribution of $X_{m+1}, \ldots, X_N$, for $|t| \leqslant CN^{1/2}$. Finally, by means of expansions we estimate the difference between the Fourier–Stieltjes transforms of $F$ and $G$.

Our proofs may be considered as an extension to the case of finite-population statistics of techniques used by Bentkus *et al.* (1997) in the i.i.d. case. We remark that the approach developed in the present paper also applies to more general nonlinear symmetric statistics based on samples drawn without replacement from finite populations. These results will appear elsewhere.

## 3.1. Notation

By $C$, $C_0$, $C_1$, $\ldots$ and $c$, $c_0$, $c_1$, $\ldots$ we denote generic absolute constants. We shall write $A \ll B$ if $A < CB$. The expression $\exp\{ix\}$ will be abbreviated to $e\{x\}$. Write

$$\Theta(t) = \left(\frac{2}{\pi}\frac{\pi - t}{\pi + t}\right)^2, \qquad \mathscr{K} = \{a \in \mathscr{A} \colon H_1|g_1(a)| < b_2\}, \qquad H_1 = \frac{b_1 N^{1/2}}{\beta_3}. \qquad (3.2)$$

Here $b_1$ is the same constant as in (1.5) and $b_2$ denotes a sufficiently small absolute constant.

Let $\nu = \{\nu_1, \ldots, \nu_n\}$ be a sequence of independent Bernoulli random variables with probabilities $P\{\nu_i = 1\} = p$ and $P\{\nu_i = 0\} = q$, for $i = 1, 2, \ldots, n$. Write

$$\beta(t) = Ee\{(\nu_1 - p)t\}, \qquad \tau = \sqrt{npq}, \qquad \delta = \delta(b_1/\beta_3, N^{1/2}),$$

where $\delta(\cdot, \cdot)$ is defined by (2.1). Let $\bar{A} = (A_1, A_2, \ldots, A_n)$ denote a random permutation which is uniformly distributed on the permutations of the ordered set $(a_1, \ldots, a_n)$ of elements of $\mathscr{A}$, independent of $\nu$. By $E^*$ we denote the conditional expectation given $\bar{A}$, i.e. $E^*(\cdot) = E(\cdot|\bar{A})$. Fo $k = 1, 2, \ldots,$ write $\Omega_k = \{1, \ldots, k\}$ and $D_k = \Omega_N \backslash \Omega_k$. Given

$D = \{i_1, i_2, \ldots, i_k\} \subset \Omega_n$, $\mathrm{E}^{i_1, \ldots, i_k}$ and $\mathrm{E}^D$ denote the conditional expectation given $A_{i_1}, \ldots, A_{i_k}$.

## 3.2 Proof of Theorem 1.1

We may and shall assume that, for sufficiently small $c_0 > 0$,

$$\frac{\beta_4}{qN} < c_0, \qquad \frac{\ln N}{\delta N} < c_0, \qquad \frac{\gamma_2}{\delta^2 q^2 N} < c_0, \qquad \frac{\ln N}{\delta qn} < c_0. \tag{3.3}$$

Indeed, if (3.3) fails, then the bound of Theorem 1.1 follows from the inequalities $F(x) \leq 1$, $|G(x)| \ll 1 + q^{-1/2}\beta_4^{1/2}/N^{1/2} + q^{1/2}\gamma_2/N^{1/2}$ and $\rho \leq \delta$; see Lemma 2.1.

   *Step 1.* Fix an integer $m \approx C_0 \delta^{-1} \ln N$, with sufficiently large $C_0$, and write

$$\Lambda_m = \sum_{1 \leq k < l \leq m} g_2(X_k, X_l). \tag{3.4}$$

Note that $U = \Lambda_m + U_1$, where $U_1$ is given by (3.1). Let $F_1$ denote the probability distribution function of $U_1$ and $\Delta_1 = \sup_x |F_1(x) - G(x)|$. We have

$$\Delta \leq \Delta_1 + \mathrm{P}\{|\Lambda_m| \geq N^{-1}\delta^{-3/2}\} + \delta^{-3/2}N^{-1}\max_x |G'(x)|.$$

By Chebyshev's inequality and the inequality $\mathrm{E}|\Lambda_m|^3 \ll m^6 \mathrm{E}|g_2(X_1, X_2)|^3$,

$$\mathrm{P}\{|\Lambda_m| \geq N^{-1}\delta^{-3/2}\} \leq \delta^{9/2} N^3 \mathrm{E}|\Lambda_m|^3 \ll \delta^{-3/2}\gamma_3 N^{-3/2} \ln^6 N.$$

Finally, using the bound $|G'(x)| \ll \beta_4/q + \gamma_2$ we obtain

$$\Delta \ll \Delta_1 + N^{-1}\delta^{-3/2}(\beta_4/q + \gamma_2 + \gamma_3). \tag{3.5}$$

Therefore, in order to prove the theorem it suffices to bound $\Delta_1$.

   Let $k$ be an integer approximately equal to $(N + m)/2$. Put $\mathcal{J}_0 = \{m + 1, \ldots, N\}$, $\mathcal{J}_0 = \Omega_n \backslash \mathcal{J}_0$, $\mathcal{J}_1 = \mathcal{J}_0 \cup \{m+1, \ldots, k\}$ and $\mathcal{J}_2 = \mathcal{J}_0 \cup \{k+1, \ldots, N\}$. Given $A$, define (random) subpopulations $\mathcal{A}_i = \{A_k, k \in \mathcal{J}_i\}$, $i = 0, 1, 2$, and let $A_i^*$ be random variables uniformly distributed in $\mathcal{A}_i$, $i = 0, 1, 2$, indpendent of $\nu$. Write

$$\upsilon_1(a) = \sum_{j=k+1}^{N} g_2(a, A_j), \qquad \upsilon_2(a) = \sum_{j=m+1}^{k} g_2(a, A_j), \tag{3.6}$$

$$H = N\delta/(32q^{-1}N(\Theta_1 + \Theta_2) + 1), \qquad \Theta_i = \mathrm{E}^*|\upsilon_i(A_i^*)|, \quad i = 1, 2.$$

Notice that $\Theta_1$ is a function of the random variables $A_{k+1}, \ldots, A_N$, and that $\Theta_2$ is a function of $A_{m+1}, \ldots, A_k$.

   *Step 2.* Split the sample as follows. Put $X_j = A_j$, for $m < j \leq N$. The rest of the sample, $X_1, \ldots, X_m$, is obtained by simple random sampling without replacement from the (random) subpopulation $\mathcal{A}_0$.

   An application of Prawitz's (1972) smoothing lemma conditionally, given $X_{m+1}, \ldots, X_m$, or equivalently, given $A_{m+1}, \ldots, A_N$, gives

$$F_1(x+) \leqslant \frac{1}{2} + \mathrm{EVP} \int_{\mathbb{R}} \mathrm{e}\{-xt\} \frac{1}{H} K\left(\frac{t}{H}\right) f_1(t)\,\mathrm{d}t, \tag{3.7}$$

$$F_1(x-) \geqslant \frac{1}{2} - \mathrm{EVP} \int_{\mathbb{R}} \mathrm{e}\{-xt\} \frac{1}{H} K\left(\frac{-t}{H}\right) f_1(t)\,\mathrm{d}t, \tag{3.8}$$

where $2K(s) = K_1(s) + \mathrm{i}K_2(s)/(\pi s)$; see, for example, Bentkus *et al.* (1997). Here

$$K_1(s) = \mathrm{I}\{|s| \leqslant 1\}(1 - |s|) \quad \text{and} \quad K_2(s) = \mathrm{I}\{|s| \leqslant 1\}(1 - |s|)\pi s \cot(\pi s).$$

Combining (3.7) and the inversion formula,

$$G(x) = \frac{1}{2} + \frac{\mathrm{i}}{2\pi} \lim_{M\to\infty} \mathrm{VP} \int_{|t|\leqslant M} \mathrm{e}\{-tx\} \hat{G}(t) \frac{\mathrm{d}t}{t}, \tag{3.9}$$

we obtain (see, for example, Bentkus *et al.* 1997)

$$F_1(x+) - G(x) \leqslant \mathrm{E}I_1 + \mathrm{E}I_2 + \mathrm{E}I_3, \tag{3.10}$$

$$I_1 = \frac{1}{2} H^{-1} \int_{\mathbb{R}} \mathrm{e}\{-xt\} K_1\left(\frac{t}{H}\right) f_1(t)\,\mathrm{d}t,$$

$$I_2 = \frac{\mathrm{i}}{2\pi} \mathrm{VP} \int_{\mathbb{R}} \mathrm{e}\{-xt\} K_2\left(\frac{t}{H}\right) (f_1(t) - \hat{G}(t)) \frac{\mathrm{d}t}{t},$$

$$I_3 = \frac{\mathrm{i}}{2\pi} \mathrm{VP} \int_{\mathbb{R}} \mathrm{e}\{-xt\} \left( K_2\left(\frac{t}{H}\right) - 1 \right) \tilde{G}(t) \frac{\mathrm{d}t}{t},$$

where VP means also that one should take $\lim_{M\to\infty}$ if necessary.

Combining (3.8) and (3.9), we obtain a bound for $G(x) - F_1(-x)$ similar to (3.10). We shall bound $F_1(x+) - G(x)$ only. To this end, we prove that

$$|\mathrm{E}I_1| + |\mathrm{E}(I_2 + I_3)| \ll N^{-1}(\beta_4/q + \delta^{-1}(\delta^{-1} + q^{-1}) + \delta^{-2}q^{-2}(\gamma_2^{1/2} + \gamma_2) + \gamma_4). \tag{3.11}$$

The analogous bound for $G(x) - F_1(x-)$ can be derived in the same way. Using these bounds, (3.5) and the inequality $\delta \geqslant \rho$ (see Lemma 2.1), we obtain the estimate of the theorem. in the remaining part of the prooof we verify (3.11).

*Step 3:* Estimate for $|\mathrm{E}I_1|$. We shall replace the random bound $H$ in the integral $I_1$ by a non-random one and $K_1(t/H)$ by 1. We have $|\mathrm{E}I_1| \leqslant |\mathrm{E}I_4| + \mathrm{E}I_5$, where

$$I_4 = H^{-1} \int_{\mathscr{Z}} \mathrm{e}\{-tx\} K_1\left(\frac{t}{H}\right) f_1(t)\,\mathrm{d}t, \qquad \mathscr{Z} = \{t \in \mathbb{R}: |t| \leqslant H_1\},$$

$$I_5 = H^{-1} \int_{H_1\leqslant|t|} K_1\left(\frac{t}{H}\right) |f_1(t)|\,\mathrm{d}t \leqslant H^{-1} \int_{H_1\leqslant|t|\leqslant H} |f_1(t)|\,\mathrm{d}t.$$

Next we construct bounds for $\mathrm{E}I_5$ and $|\mathrm{E}I_4|$; see (3.12) and (3.19) below. It follows from these bounds that $|\mathrm{E}I_1|$ does not exceed the right-hand side of (3.11).

Let us show

$$\mathrm{E}I_5 \ll N^{-2}\beta_3. \qquad (3.12)$$

For this purpose we represent $f_1(t)$ in Erdős and Rényi (1959) form conditionally, given $A_{m+1}, \ldots, A_N$. Let $v^* = \{v_1^*, \ldots, v_n^*\}$ be a sequence of independent Bernoulli random variables independent of $\mathscr{A}$ and with probabilities

$$\mathrm{P}\{v_i^* = 1\} = p^*, \qquad \mathrm{P}\{v_i^* = 0\} = q^*, \qquad p^* = \frac{m}{n - (N - m)}, \qquad q^* = 1 - p^*.$$

Write $S_* = \sum_{k \in \mathscr{J}_0}(v_k^* - p^*)$ and $L_* = \sum_{k \in \mathscr{J}_0} l(A_k)v_k^*$. We have

$$f_1(t) = \mathrm{P}^{-1}\{S_* = 0\}\frac{1}{2\pi}\int_{-\pi}^{\pi} W \, \mathrm{d}s, \qquad W = \mathrm{E}^*\mathrm{e}\{t(L_* + U') + sS_*\}. \qquad (3.13)$$

We shall construct an upper bound for $|W|$. We have

$$|W| = \prod_{k \in \mathscr{J}_0} |\beta_*(z(A_k) + tv(A_k))|, \qquad \beta_*(x) = \mathrm{E}\mathrm{e}\{(v_1^* - p^*)x\}.$$

Here we denote

$$z(a) = tg_1(a) + s \quad \text{and} \quad v(a) = v_1(a) + v_2(a),$$

with $v_i(a)$ given by (3.6). Then we apply the identity $|\beta_*(x)|^2 = 1 - 2p^*q^*(1 - \cos x)$ to $x = z(a) + tv(a)$ and expand the cosine function in powers of $tv(a)$ to obtain

$$|\beta_*(z(a) + tv(a))|^2 \leqslant u_1(a) + u_2(a), \qquad (3.14)$$

$$u_1(a) = 1 - 2p^*q^*(1 - \cos(z(a))), \qquad u_2(a) = 2p^*q^*|tv(a)|.$$

Furthermore, we may assume that $p^* \leqslant 8^{-1}$ (this is a consequence of the last inequality of (3.3) provided that $c_0$ is small enough). This inequality implies $u_1(a) \geqslant 1/2$ and, therefore,

$$u_1(a) + u_2(a) \leqslant u_1(a)(1 + 2u_2(a)). \qquad (3.15)$$

Combining (3.14) and (3.15), we obtain

$$|W|^2 \leqslant W_1 W_2, \qquad W_1 = \prod_{k \in \mathscr{J}_0} u_1(A_k), \qquad W_2 = \prod_{k \in \mathscr{J}_0}(1 + 2u_2(A_k)). \qquad (3.16)$$

To estimate $W_2$, we apply the arithmetic-geometric mean inequality,

$$W_2 \leqslant \left(\frac{1}{|\mathscr{J}_0|}\sum_{k \in \mathscr{J}_0}(1 + 2u_2(A_k))\right)^{|\mathscr{J}_0|} = (\mathrm{E}^*(1 + 2u_2(A_0^*)))^{n-N+m}, \qquad (3.17)$$

and use (5.2) to bound $\mathrm{E}^*|v(A_0^*)| \leqslant q^{-1}(\Theta_1 + \Theta_2)$. Thus, for $|t| \leqslant H$, we obtain

$$\mathrm{E}^*(1 + 2u_2(A_0^*)) \leqslant 1 + 4p^*q^*q^{-1}H(\Theta_1 + \Theta_2) \leqslant 1 + p^*q^*\frac{\delta}{8} \leqslant \exp\left\{p^*q^*\frac{\delta}{8}\right\}.$$

This inequality, in combination with (3.17), implies $W_2^{1/2} \leqslant \exp\{mq^*\delta/16\}$. Now in view of (3.16) and (3.13) we obtain, for $|t| \leqslant H$,

$$|f_1(t)| \ll W_3 W_1^{1/2}, \qquad W_3 = m^{1/2}\exp\{mq^*\delta/16\}.$$

Here we have estimated $P^{-1}\{S_* = 0\} \ll m^{1/2}$; see (5.16). We have

$$\mathrm{E}I_5 \leqslant \frac{1}{H_1}\mathrm{E}\int_{H_1 \leqslant |t| \leqslant H}|f_1(t)|\,\mathrm{d}t \leqslant \frac{W_3}{H_1}\int_{H_1 \leqslant |t| \leqslant N}\mathrm{E}W_1^{1/2}\,\mathrm{d}t. \tag{3.18}$$

To bound $\mathrm{E}W_1^{1/2}$ we apply Hölder's inequality and Theorem 4 of Hoeffding (1963),

$$(\mathrm{E}W_1^{1/2})^2 \leqslant \mathrm{E}W_1 \leqslant (\mathrm{E}u_1(A_1))^{|\mathscr{I}_0|};$$

see Section 5 below. Note that $\mathrm{E}u_1(A_1) \leqslant 1 - 2p^*q^*\delta$, for $H_1 \leqslant |t| \leqslant N$, by the choice of $\delta$. Therefore,

$$\mathrm{E}W_1^{1/2} \leqslant (1 - 2p^*q^*\delta)^{(n-N+m)/2} \leqslant \exp\{-p^*q^*\delta(n-N+m)\} = \exp\{-mq^*\delta\}.$$

Combining this bound with (3.18) and using the inequality $q^* = 1 - p^* \geqslant 7/8$, we obtain (3.12), provided that the constant $C_0$ (in the definition of $m$) is sufficiently large.

We must now bound $\mathrm{E}I_4$. We shall show

$$|\mathrm{E}I_4| \ll \mathscr{R}_0, \qquad \mathscr{R}_0 = N^{-1}\delta^{-2}(1 + q^{-2}\gamma_2) + N^{-1}\delta^{-1}q^{-1}(1 + q^{-1}\gamma_2^{1/2}). \tag{3.19}$$

It follows from the inequality $|K_1(u) - 1| \leqslant |u|$ that

$$I_4 = I_6 + R, \qquad I_6 = H^{-1}\int_{\mathscr{Z}}\mathrm{e}\{-tx\}f_1(t)\,\mathrm{d}t, \tag{3.20}$$

$$\mathrm{E}|R| \leqslant \mathrm{E}H^{-1}\int_{\mathscr{Z}}|t|H^{-1}\,\mathrm{d}t = H_1^2\mathrm{E}H^{-2} \ll \mathscr{R}_0,$$

where in the last step we have applied (5.1). Recall that $U = U_1 + \Lambda_m$. Now, using the inequality $|\mathrm{e}\{t\Lambda_m\} - 1| \leqslant |t\Lambda_m|$, we obtain

$$I_6 = I_7 + R, \qquad I_7 = H^{-1}\int_{\mathscr{Z}}\mathrm{e}\{-tx\}f_2(t)\,\mathrm{d}t, \qquad f_2(t) = \mathrm{E}^{D_m}\mathrm{e}\{tU\}, \tag{3.21}$$

$$\mathrm{E}|R| \leqslant \mathrm{E}H^{-1}\int_{\mathscr{Z}}\mathrm{E}^{D_m}|t\Lambda_m|\,\mathrm{d}t \leqslant H_1^2\mathrm{E}H^{-1}|\Lambda_m| \ll \mathscr{R}_0,$$

where in the last step we have used the inequality $|\Lambda_m H^{-1}| \leqslant \Lambda_m^2 + H^{-2}$ and moment inequalities (5.1) and (5.3). Next we replace $I_7$ by

$$I_8 = H^{-1}\int_{\mathscr{Z}_0}\mathrm{e}\{-tx\}f_2(t)\,\mathrm{d}t, \qquad \mathscr{Z}_0 = \{C_1q^{-1} \leqslant |t| \leqslant H_1\}, \tag{3.22}$$

where $C_1$ is a sufficiently large constant. We have $I_7 = I_8 + R$ with $|R| \leqslant 2C_1q^{-1}H^{-1}$. Hölder's inequality, in combination with (5.1), gives $\mathrm{E}|R| \ll \mathscr{R}_0$.

It remains to estimate $\mathrm{E}I_8$. Write $I_8 = 32q^{-1}\delta^{-1}(J_1 + J_2) + \delta^{-1}J_3$, where

$$J_i = \int_{\mathscr{Z}_0}\mathrm{e}\{-tx\}f_2(t)\Theta_i\,\mathrm{d}t, \quad i = 1, 2, \qquad J_3 = N^{-1}\int_{\mathscr{Z}_0}\mathrm{e}\{-tx\}f_2(t)\,\mathrm{d}t.$$

In order to complete the prooof of (3.19), we shall show

$$\mathrm{E}J_i \ll N^{-1}(1 + \gamma_2), \qquad i = 1, 2, 3. \tag{3.23}$$

Let us prove (3.23) for $i = 1, 2$. By symmetry, it suffices to consider the case where $i = 1$. Recall that the random variable $\Theta_1$ is a function of $X_{k+1}, \ldots, X_N$. In view of the inequality $k \approx (N + m)/2 > m$, we can write

$$\mathrm{E}\Theta_1 f_2(t) = \mathrm{E}\Theta_1 f_3(t), \qquad f_3(t) = \mathrm{E}(\mathrm{e}\{tU\}|X_{k+1}, \ldots, X_N). \tag{3.24}$$

Given $t \in \mathscr{Z}_0$, choose an integer $m_1 = C_2 N t^{-2} \ln|t|$. Here $C_2$ is a sufficiently large constant to be specified later. Given $C_2$, we may choose $C_1$ in (3.22) large enough so that $m_1 < 10^{-1} qN < k$, for $t \in \mathscr{Z}_0$. Write $\mathscr{J}_3 = \Omega_{m_1} \cap (\Omega_n \setminus \Omega_N)$. We shall represent our sample $X_1, \ldots, X_N$ as follows. For $m_1 + 1 \leq j \leq N$, put $X_j = A_j$. The remaining part of the sample (the observations $X_1, \ldots, X_{m_1}$) represents a simple random sample drawn without replacement form the set $\mathscr{A}_3 = \{A_k, k \in \mathscr{J}_3\}$. Let $A_3^*$ be a random variable uniformly distributed in $\mathscr{A}_3$. Put

$$v_3(a) = \sum_{k=m_1+1}^{N} g_2(a, A_j) \quad \text{and} \quad \Theta_3 = \mathrm{E}^* |v_3(A_3^*)|. \tag{3.25}$$

Notice that the random variable $\Theta_3$ is a function of $A_{m_1+1}, \ldots, A_N$.

Write $U = U_1^\star + \Lambda_{m_1}$, where $U_1^\star = L'_\star + U'_\star$, with

$$L'_\star = l_\star(X_1) + \ldots + l_\star(X_{m_1}), \qquad l_\star(x) = g_1(x) + l_0^\star(x), \qquad l_0^\star(x) = \sum_{j=m_1+1}^{N} g_2(x, X_j),$$

and with $U'_\star$ defined by (3.1), but with $m$ replaced by $m_1$. Furthermore, $\Lambda_{m_1}$ is given by (3.4).

Using the inequality $|\mathrm{e}\{t\Lambda_{m_1}\} - 1| \leq |t\Lambda_{m_1}|$, we obtain

$$\mathrm{E}\Theta_1 f_3(t) = \mathrm{E}\Theta_1 f_4(t) + R_1, \qquad f_4(t) = \mathrm{E}(\mathrm{e}\{tU_1^*\}|X_{k+1}, \ldots, X_N), \tag{3.26}$$

where $|R_1| \leq \mathrm{E}|t\Lambda_{m_1}|\Theta_1$. Furthermore, combining (3.24) and (3.26), we obtain

$$\mathrm{E}J_1 = \mathrm{E}J_4 + R, \qquad J_4 = \int_{\mathscr{Z}_0} \mathrm{e}\{-tx\} f_4(t)\Theta_1 \, \mathrm{d}t, \tag{3.27}$$

$$|R| \leq \int_{\mathscr{Z}_0} \mathrm{E}|t\Lambda_{m_1}|\Theta_1 \, \mathrm{d}t \ll N^{-1}\gamma_2.$$

In the last step we invoke (5.1), (5.3) and apply Hölder's inequality to obtain

$$\mathrm{E}|\Lambda_{m_1}|\Theta_1 \leq (\mathrm{E}\Lambda_{m_1}^2)^{1/2}(\mathrm{E}\Theta_1^2)^{1/2} \ll m_1 N^{-5/2}\gamma_2 \ll t^{-1}\ln|t| N^{-3/2}\gamma_2$$

and bound the integral of the function $|t|^{-1}\ln|t|$ over the region $\mathscr{Z}_0$ by $\ln^2 N$.

To estimate $\mathrm{E}J_4$ observe that, by the inequality $m_1 < k$,

$$\mathrm{E}\Theta_1 f_4 = \mathrm{E}\Theta_1 f_5, \qquad f_5 = \mathrm{E}(\mathrm{e}\{tU'_\star\}|X_{m_1+1}, \ldots, X_N).$$

Therefore, $\mathrm{E}J_4 = \mathrm{E}J_5$, where $J_5$ is defined in the same way as $J_4$ (see (3.27)) but with $f_4$ replaced by $f_5$. Furthermore,

$$\mathrm{E}J_5 = \mathrm{E}J_6 + R, \qquad J_6 = \int_{\mathscr{Z}_0} \mathrm{e}\{-tx\} f_5(t)\Theta_1 I_\Theta \, \mathrm{d}t, \tag{3.28}$$

$$I_\Theta = I\{N\Theta_3 \leqslant c_1|t|\}, \qquad |R| \leqslant \int_{\mathscr{Z}_0} \mathrm{E}\Theta_1 I\{N\Theta_3 > c_1|t|\}\,\mathrm{d}t \ll N\mathrm{E}\Theta_1\Theta_3.$$

Here $c_1$ denotes a small positive constant to be determined below. Combining (5.3) and Hölder's inequality, we obtain $|R| \ll N^{-1}\gamma_2$.

In order to bound $\mathrm{E}J_6$ we represent $f_5$ in the Erdős and Rényi (1959) form; see (3.29). Let $\nu^\star = \{\nu_1^\star, \ldots, \nu_n^\star\}$ be a sequence of independent Bernoulli random variables independent of $\mathscr{A}$ and with probabilities

$$\mathrm{P}\{\nu_i^\star = 1\} = p^\star, \qquad \mathrm{P}\{\nu_i^\star = 0\} = q^\star, \qquad p^\star = \frac{m_1}{n - (N - m_1)}, \qquad q^\star = 1 - p^\star.$$

Write $S_\star = \sum_{k \in \mathscr{J}_3}(\nu_k^\star - p^\star)$, $L_\star = \sum_{k \in \mathscr{J}_3} l_\star(A_k)\nu_k^\star$ and $\tau_\star^2 = m_1 q^\star$. We have

$$f_5(t) = \lambda_\star \int_{-\pi\tau_\star}^{\pi\tau_\star} W_\star\,\mathrm{d}s, \qquad W_\star = \mathrm{E}^\star \mathrm{e}\left\{t(L_\star + U_\star') + \frac{s}{\tau_\star}S_\star\right\}, \tag{3.29}$$

with $\lambda_\star^{-1} = 2\pi\tau_\star \mathrm{P}\{S_\star = 0\}$ satisfying $\lambda_\star \ll 1$, by (5.16).

Combining (3.28) and (3.29), we obtain

$$\mathrm{E}J_6 \ll \int_{\mathscr{Z}_0}\mathrm{d}t \int_{-\pi\tau_\star}^{\pi\tau_\star} \mathrm{E}\Theta_1 I_\Theta |W_\star|\,\mathrm{d}s. \tag{3.30}$$

In the next step we construct an upper bound for $\mathrm{E}\Theta_1 I_\Theta |W_\star|$. Note that the inequality $m_1 < 10^{-1}qN$ implies $p^\star \leqslant 10^{-1}$. The same argument as above (see (3.16)) gives

$$|W_\star|^2 \leqslant W_1^\star W_2^\star, \qquad W_1^\star = \prod_{k \in \mathscr{J}_3} u_1^\star(A_k), \qquad W_2^\star = \prod_{k \in \mathscr{J}_3}(1 + 2u_2^\star(A_k)), \tag{3.31}$$

where $u_1^\star$ and $u_2^\star$ are given by (3.14), but with $p^*$, $q^*$, $z(a)$ and $v(a)$ replaced by $p^\star$, $q^\star$, $z_\star(a) := tg_1(a) + s/\tau_\star$ and $v_3(a)$ (defined in (3.25)) respectively.

To bound $W_2^\star$ we proceed as in (3.17) and obtain

$$W_2^\star \leqslant (1 + 2\mathrm{E}^\star u_2^\star(A_3^*))^{n-N+m_1} = (1 + 4p^\star q^\star|t|\Theta_3)^{n-N+m_1} \leqslant \exp\{4m_1 q^\star|t|\Theta_3\}.$$

Furthermore, by our choice of $m_1$, $I_\Theta(W_2^\star)^{1/2} \leqslant \exp\{2q^\star C_2 c_1 \ln|t|\}$. Therefore, in view of (3.31),

$$\mathrm{E}\Theta_1 I_\Theta |W_\star| \leqslant \exp\{2q^\star C_2 c_1 \ln|t|\}\mathrm{E}\Theta_1(W_1^\star)^{1/2}. \tag{3.32}$$

Now we apply Hölder's inequality and invoke (5.1) to obtain

$$\mathrm{E}\Theta_1(W_1^\star)^{1/2} \leqslant (\mathrm{E}\Theta_1^2)^{1/2}(\mathrm{E}W_1^\star)^{1/2} \ll N^{-1}\gamma_2^{1/2}(\mathrm{E}W_1^\star)^{1/2}. \tag{3.33}$$

To bound $\mathrm{E}W_1^\star$ we apply Theorem 4 of Hoeffding (1963) and obtain

$$\mathrm{E}W_1^\star \leqslant (\mathrm{E}u_1^\star(A_1))^{|\mathscr{J}_3|} = (1 - 2p^\star q^\star M)^{n-N+m_1} \leqslant \exp\{-2m_1 q^\star M\}, \tag{3.34}$$

where $M = \mathrm{E}(1 - \cos z_\star(A_1))$. Combining the inequalities

$$M \geqslant \mathrm{E}(1 - \cos z_\star(A_1))I_{\mathscr{K}}(A_1), \qquad I_{\mathscr{K}}(a) = I\{a \in \mathscr{K}\},$$

$$1 - \cos z_\star(a) \geqslant 2^{-1}\Theta(b_2)z_\star^2(a), \qquad a \in \mathscr{K},$$

(see (5.15)) we get $M \geqslant 2^{-1}\Theta(b_2)\mathrm{E}z_\star^2(A_1)I_{\mathscr{H}}(A_1)$. Now, by Lemma 5.3,

$$M \geqslant b_3(t^2 N^{-1} + s^2 \tau_\star^{-2}), \qquad b_3 = 2^{-1}\Theta(b_2)(1 - 2b_1 b_2^{-1}),$$

is a positive constant (because of our choice of $0 < 2b_1 < b_2$ in (1.5) and (3.2)). Substituting this inequality in (3.34) and using $q^\star = 1 - p^\star \geqslant 9/10$, we obtain

$$\mathrm{E}W_1^\star \leqslant \exp\{-2b_3 m_1 q^\star(t^2 N^{-1} + s^2 \tau_\star^{-2})\} \leqslant \exp\{-2b_3(\tfrac{9}{10}C_2 \ln|t| + s^2)\}. \tag{3.35}$$

Finally, collecting the inequalities (3.32), (3.33) and (3.35) in (3.30), we obtain

$$\mathrm{E}J_6 \ll N^{-1}\gamma_2^{1/2}\int_{\mathscr{Z}_0} \mathrm{d}t \int_{-\pi\tau_\star}^{\pi\tau_\star} \exp\{C_2(2c_1 - \tfrac{9}{10}b_3)\ln|t| - b_3 s^2\}\,\mathrm{d}s. \tag{3.36}$$

Choosing $c_1 = b_3/4$ and $C_2 = 4/b_3$, we obtain bounded integrals in (3.36), and thus $\mathrm{E}J_6 \ll N^{-1}\gamma_2^{1/2} \leqslant N^{-1}(1 + \gamma_2)$. This inequality, together with (3.27) and (3.28), completes the proof of (3.23) in the case where $i = 1$.

The proof of (3.23) in the case where $i = 3$ is similar but simpler: just write $N^{-1}$ instead of $\Theta_1$ in the proof above.

Collecting the bounds (3.20), (3.21), (3.22) and (3.23), we obtain (3.19).

*Step 4.* Estimate for $|\mathrm{E}(I_2 + I_3)|$. Write $I_2 + I_3 = \mathrm{i}(2\pi)^{-1}(I_9 + I_{10} - I_{11} + I_{12})$, where

$$I_9 = \int_{|t| \leqslant H_1} \mathrm{e}\{-tx\}\frac{f_1(t) - \hat{G}(t)}{t}\,\mathrm{d}t, \quad I_{10} = \int_{H_1 \leqslant |t| \leqslant H} \mathrm{e}\{-tx\}K_2\left(\frac{t}{H}\right)f_1(t)\frac{\mathrm{d}t}{t},$$

$$I_{11} = \int_{|t| > H_1} \mathrm{e}\{-tx\}\hat{G}(t)\frac{\mathrm{d}t}{t}, \quad I_{12} = \int_{|t| \leqslant H_1} \mathrm{e}\{-tx\}\left(K_2\left(\frac{t}{H}\right) - 1\right)f_1(t)\frac{\mathrm{d}t}{t}.$$

Using (3.3), it is easy to show that $|\mathrm{E}I_{11}| \ll q^{-1}\beta_4/N + \gamma_2/N$. Using the inequality $|K_2(s) - 1| \leqslant cs^2$, and invoking (5.1), we obtain

$$|\mathrm{E}I_{12}| \ll \mathrm{E}H^{-2}H_1^2 \ll \delta^{-2}N^{-1}(1 + q^{-2}\gamma_2).$$

To bound $|\mathrm{E}I_{10}|$ write

$$|\mathrm{E}I_{10}| \leqslant \mathrm{E}I_{13}, \quad I_{13} = \int_{H_1 \leqslant |t| \leqslant H} |f_1(t)|\frac{\mathrm{d}t}{|t|}.$$

The bound $\mathrm{E}I_{13} \ll N^{-1}\beta_3$ is obtained in a similar way as (3.12) above. Collecting these inequalities, we obtain

$$|\mathrm{E}(I_2 + I_3)| \ll |\mathrm{E}I_9| + N^{-1}q^{-1}\beta_4 + N^{-1}\delta^{-2}(1 + q^{-2}\gamma_2). \tag{3.37}$$

In order to complete the proof of (3.11) we shall show that

$$|\mathrm{E}I_9| \ll \delta^{-2}N^{-1}(1 + \gamma_2) + N^{-1}(q^{-1}\beta_4 + \gamma_4). \tag{3.38}$$

We have

$$\mathrm{E}I_9 = \int_{|t| \leqslant H_1} \mathrm{e}\{-tx\}(\mathrm{E}\mathrm{e}\{tU_1\} - \hat{G}(t))\frac{\mathrm{d}t}{t}.$$

Recall that $U_1 = U - \Lambda_m$. Write $e\{tU_1\} = e\{tU\}e\{-t\Lambda_m\}$ and expand $e\{-t\Lambda_m\}$ in powers of $-it\Lambda_m$ to obtain $EI_9 = I_{14} - iI_{15} + R$, where

$$I_{14} = \int_{-H_1}^{H_1} e\{-tx\}\frac{\hat{F}(t) - \hat{G}(t)}{t}\,dt, \quad I_{15} = \int_{-H_1}^{H_1} e\{-tx\}E\Lambda_m e\{tU\}\,dt,$$

and where $|R| \leq H_1^2 E\Lambda_m^2 \ll \delta^{-2} N^{-1}\gamma_2$, by (5.3). By symmetry, $EI_{15} = \binom{m}{2}EI_{16}$, where $I_{16}$ is defined in the same way as $I_{15}$, but with $\Lambda_m$ replaced by $g_2(X_{N-1}, X_N)$. the bound $EI_{16} \ll N^{-3/2}(1 + \gamma_2)$ is obtained in a similar way to (3.23): just take $\tilde{f}_3 = E(e\{tU\}|X_{N-1}, X_N)$ instead of $f_3$ and $g_2(X_{N-1}, X_N)$ instead of $\Theta_1$ in the proof of (3.23) (for $i = 1$). We obtain

$$|EI_9 - I_{14}| \ll \delta^{-2}N^{-1}(1 + \gamma_2).$$

In the next section on expansions (see (4.1) below) we shall show $|I_{14}| \ll N^{-1}(\beta_4/q + \gamma_4)$, thus completing the proof of (3.38).

## 3.3. Proof of Theorem 1.2

The bound of the theorem follows from (3.11). Just note that for a linear statistic we have $g_2(x, y) = 0$, for any $x, y \in \mathscr{A}$. In particular, we do not need to assume that the last two inequalities of (3.3) hold.

## 3.4. Proof of Corollary 1.3

The corollary follows from Theorem 1.1, by the law of large numbers (LLN) for $U$-statistics; see, for example, Serfling (1980). Given $N$, the function $\mathscr{H}$ and a sequence of i.i.d. observations $\mathscr{X}_1, \mathscr{X}_2, \ldots$, introduce the sequence of finite populations $\mathscr{A}_n = \{\mathscr{X}_1, \ldots, \mathscr{X}_n\}$ and the corresponding sequence of $U$-statistics, $(U_n)$. Given $x \in \mathbb{R}$, apply the bound of Theorem 1.1 to the sequence of probabilities $P_n\{x\} = P\{U_n \leq x\}$. By the LLN, we obtain $\lim_n P_n\{x\} = P\{\tilde{U} \leq x\}$. Furthermore, the moments of the linear and quadratic parts of $U_n$ in the expansion and in the remainder (in the estimate of Theorem 1.1) converge to the corresponding moments of the statistics $\tilde{U}$, thus proving Corollary 1.3.

# 4. Expansions

As in the previous section, we assume that $\beta_2 = 1$ and that inequalities (3.3) hold. With $H_1$ given in (3.2), in this section we shall prove the inequality

$$\int_{|t| \leq H_1} |t|^{-1}|\hat{F}(t) - \hat{G}(t)|\,dt \ll \mathscr{R}, \qquad \mathscr{R} := \frac{1}{N}\left(\frac{\beta_4}{q} + \gamma_4\right). \tag{4.1}$$

We introduce some notation. Let $\theta_1, \theta_2, \ldots$ denote independent random variables

uniformly distributed in [0, 1] and independent of all other random variables considered. For a vector-valued smooth function $H$ we use the Taylor expansion

$$H(x) = H(0) + H'(0)x + \ldots + H^{(n)}(0)\frac{x^n}{n!} + \mathrm{E}_{\theta_1} H^{(n+1)}(\theta_1 x)(1 - \theta_1)^n \frac{x^{n+1}}{n!}.$$

Here $\mathrm{E}_{\theta_1}$ denotes the conditional expectation given all the random variables but $\theta_1$. In particular, we have the mean value formula, $H(x) - H(0) = \mathrm{E}_{\theta_1} H'(\theta_1 x)x$.

Given a sum $S = s_1 + \ldots + s_k$, denote $S^{(i)} = S - s_i$ and, similarly, $S^{(i,j)} = S - s_i - s_j$.

Using the fact that the distribution of $U$ coincides with the conditional distribution of

$$U_0 := \sum_{1 \leq i < j \leq n} h(A_i, A_j)\nu_i\nu_j$$

$$= \sum_{i=1}^{n} g_1(A_i)(\nu_i - p) + \sum_{1 \leq i < j \leq n} g_2(A_i, A_j)(\nu_i - p)(\nu_j - p),$$

conditioned on the event $\mathbb{B} := \{S_0 = N\}$, where $S_0 = \sum_{i=1}^{n} \nu_i$, we obtain

$$\hat{F}(t) = \frac{1}{2\pi \mathrm{P}\{\mathbb{B}\}} \int_{-\pi}^{\pi} \mathrm{E}e\{tU_0 + s(S_0 - N)\} \, \mathrm{d}s;$$

see Erdős and Rényi (1959). Write

$$T = \sum_{i=1}^{n} T_i, \qquad T_i = z_i(\nu_i - p), \qquad z_i = tx_i + s\tau^{-1}, \qquad x_i = g_1(A_i), \qquad \tau = (npq)^{1/2},$$

$$Q = \sum_{1 \leq i < j \leq n} Q_{i,j}, \qquad Q_{i,j} = ty_{i,j}(\nu_i - p)(\nu_j - p), \qquad y_{i,j} = g_2(A_i, A_j).$$

We have $T + Q = tU_0 + s\tau^{-1}(S_0 - N)$ and, therefore,

$$\hat{F}(t) = \lambda \int_{-\pi\tau}^{\pi\tau} \mathrm{E}e\{T + Q\} \, \mathrm{d}s, \qquad \lambda^{-1} = 2\pi\tau \mathrm{P}\{\mathbb{B}\}.$$

Höglund (1978) showed that $2^{-1/2}\pi \leq \lambda^{-1} \leq (2\pi)^{1/2}$; see (5.16). We shall approximate the integrand $\mathrm{E}e\{T + Q\}$ by the sum $h_1 + h_2$, where

$$h_1 = \mathrm{E}e\{t\}, \qquad h_2 = \mathrm{i}^3 \binom{n}{2} \mathrm{E}e\{T^{(1,2)}\}V, \qquad V = Q_{1,2}T_1T_2.$$

To prove (4.1) it clearly suffices to prove the inequalities

$$\int_{|t| \leq H_1} \left| \lambda \int_{|s| \leq \pi\tau} (h_1 + h_2) \, \mathrm{d}s - \hat{G}(t) \right| \frac{\mathrm{d}t}{|t|} \ll \mathscr{R}, \tag{4.2}$$

$$I := \int_{|t| \leq H_1} \lambda \int_{|s| \leq \pi\tau} |\mathrm{E}e\{T + Q\} - (h_1 + h_2)| \, \mathrm{d}s \frac{\mathrm{d}t}{|t|} \ll \mathscr{R}. \tag{4.3}$$

Note that in the i.i.d. case the inequality corresponding to (4.2) is proved in Lemma 6.1 of Bentkus *et al.* (1997). We prove (4.2) by combining the proof of this lemma with the proof of

the Berry–Esseen bound for the finite-population sample mean given in Höglund (1978). For details we refer to Lemma 4.3 of Bloznelis and Götze (1997).

To prove (4.3), we expand $e\{T + Q\}$ in powers of $T_i$ and $Q_{i,j}$. In order to ensure the integrability (with respect to the measure $ds\,dt/|t|$) of the remainders of these expansions we split $Ee\{T + Q\}$ into a product of two functions (different for different values of $s$ and $t$) so that the first one is the characteristic function of a sum of conditionally independent random variables and vanishes sufficiently fast as $s$ and $t$ tend to infinity. This type of approach has been used earlier by Helmers and van Zwet (1982), van Zwet (1984), Götze and van Zwet (1991) and Bentkus *et al.* (1997) in the i.i.d. situation.

Introduce the set $\mathscr{Z} = \{(s, t)\colon |s| \leqslant \pi\tau, |t| \leqslant H_1\}$. For technical reasons it is convenient to split the integral $I$ in two parts $I = I_1 I_2$ according to the regions $\mathscr{Z} = \mathscr{Z}_1 \cup \mathscr{Z}_2$,

$$Z_1 = \mathscr{Z} \cap \{|t| \leqslant C_3 q^{-1}\} \quad \text{and} \quad Z_2 = \mathscr{Z} \cap \{C_3 q^{-1} < |t| \leqslant H_1\}. \tag{4.4}$$

Here $C_3$ denotes a sufficiently large absolute constant. We choose $C_3 = 600\Theta^{-1}(1)$. In Lemma 4.1 we prove the bound $I_2 \ll \mathscr{R}$. The proof of the bound $I_2 \ll \mathscr{R}$ is similar but simpler. We skip it and refer to Lemma 4.2 Bloznelis and Götze (1997) for details. It remains to prove Lemma 4.1.

Note that, for any $i, j, i_1, \ldots, i_k \in \Omega_n$ such that $\{i, j\} \cap \{i_1, \ldots, i_k\} = \varnothing$ we have

$$\mathrm{E}^{i_1, \ldots, i_k}|y_{i,j}|^r \leqslant c(k, r)\mathrm{E}|y_{i,j}|^r, \quad \mathrm{E}^{i_1, \ldots, i_k}|x_j|^r \leqslant c(k, r)\mathrm{E}|x_j|^r, \qquad r \geqslant 0. \tag{4.5}$$

We need to introduce some more notation. Given $D = \{i, j, \ldots, k\} \subset \Omega_n$, let $\mathrm{E}_{\{D\}} = \mathrm{E}_{\{i,j,\ldots,k\}}$ and $\mathrm{E}_{[D]} = \mathrm{E}_{[i,j,\ldots,k]}$ denote the conditional expectation given all the random variables but $\{\nu_j, j \in D\}$ and the conditional expectation given $\{\nu_j, A_j, j \in D\}$, respectively. Given $1 \leqslant m \leqslant n$, introduce the random variables

$$\xi_i = t(\nu_i - p)\zeta_m(A_i), \qquad \zeta_m(a) = \sum_{j=m+1}^n g_2(a, A_j)(\nu_j - p). \tag{4.6}$$

Here $i \in \Omega_m$ and $a \in \mathscr{A} \backslash \{A_{m+1}, \ldots, A_n\}$. Given $B \subset \Omega_m$, denote

$$Y_B = \left| \mathrm{E}_{\{B\}} e\left\{\sum_{i \in B} T_i\right\} \right|, \quad Z_B = \left| \mathrm{E}_{\{B\}} e\left\{\sum_{i \in B} (T_i + \xi_i)\right\} \right|.$$

Furthermore, given $A_i, i \in B$, let $A_B^*$ denote the random variable uniformly distributed in the set $\{A_i, i \in B\}$ and let $\mathrm{E}_B^*$ denote the conditional expectation given all the random variables but $A_B^*$. Introduce the random variables

$$\Psi_B = g_B(t) \prod_{k \in B} u_{[1]}^{1/2}(z_k), \qquad \kappa_B = \alpha N \mathrm{E}_B^* \zeta_m^2(A_B^*), \qquad I_B = I\{\kappa_B > \delta\}, \tag{4.7}$$

where $\alpha = 2\pi(4\Theta^{-1}(1) + 1)$ and $\delta = \Theta(1)/40$ are constants,

$$g_B(t) = \exp\left\{pq\frac{\delta}{2}\frac{|B|}{N}t^2\right\}, \qquad u_{[d]}(x) = 1 - \frac{pq}{2}\Theta(d)x^2 I\{|x| < d + \pi\}, \quad d > 0. \tag{4.8}$$

In Lemma 5.4 below, for $|t| \leqslant H_1$ and $|s| \leqslant \pi\tau$, we prove the inequalities

$$Z_B \ll I_B + \Psi_B, \qquad Y_B \ll \Psi_B, \qquad \mathrm{E}^{i_1,\ldots,i_4}\Psi_B^r \ll F_B^r, \qquad r = 1, 2, \tag{4.9}$$

where $i_1, \ldots, i_4 \in \Omega_n \backslash B$. Here we denote

$$F_B = \exp\{-8\delta pq|B|N^{-1}(t^2 + s^2/q)\}.$$

We often take $|B| \geqslant m/4$, with $m$ given by (4.13). In this case we have

$$F_B \leqslant (t^2 + s^2/q)^{-10}. \tag{4.10}$$

**Lemma 4.1.** *Assume that* $\beta_2 = 1$ *and that (3.3) holds. Then*

$$I_2 = \lambda \int_{\mathscr{Z}_2} \frac{|\mathrm{Ee}\{T + Q\} - (h_1 + h_2)|}{|t|} \, ds \, dt \ll \mathscr{R}, \tag{4.11}$$

*where* $\mathscr{Z}_2$ *is given by (4.4).*

***Proof.*** Given a positive number $L$ and a complex-valued function $f(s, t)$, we write $f \prec L$ if

$$\int_{\mathscr{Z}_2} |f(s, t)||t|^{-1} \, ds \, dt \ll L.$$

Furthermore, for two complex-valued functions $f$, $g$ we write $f \sim g$ if $f - g \prec \mathscr{R}$. In view of the inequality $\lambda \leqslant 2^{1/2}\pi^{-1}$, (4.11) can be abbreviated as follows:

$$\mathrm{Ee}\{T + Q\} \sim h_1 + h_2. \tag{4.12}$$

Given $(s, t) \in \mathscr{Z}_2$ wirte $u = t^2 + s^2/q$ and let

$$m = m(s, t) > C_4 q^{-1} n u^{-1} \ln u, \qquad C_4 = 300\Theta^{-1}(1), \tag{4.13}$$

denote the smallest integer which is greater than $C_4 q^{-1} n u^{-1} \ln u$. A simple calculation shows that $C_4 \leqslant m(s, t) \leqslant C_4 C_3^{-1} n$, for $(s, t) \in \mathscr{Z}_2$. Since $C_4 = C_3/2$ we have $10 \leqslant m(s, t) \leqslant n/2$.

Write $\mu := mpqN^{-1} = C_4 u^{-1} \ln u$. We shall often use the following fact,

$$(t^2)^\alpha (s^2)^\beta \mu^\gamma \prec q^{\beta+1/2} c(\alpha, \beta, \gamma), \qquad \text{for } \gamma > \alpha + \beta + \tfrac{1}{2}, \qquad \alpha, \beta \geqslant 0.$$

In what follows $B$ always denote the set $\{4, \ldots, m\}$. $R, R_1, R_2 \ldots$ will denote random variables (remainders) which may be different in different places. This will not cause any misunderstanding if we assume that $R, R_1, R_2 \ldots$ always take the latest prescribed values.

Let us prove (4.12). Split $Q = Q_A + Q_D + \xi$ and $T = T_A + T_D$, where

$$Q_A = \sum_{1 \leqslant i < j \leqslant m} Q_{i,j}, \qquad Q_D = \sum_{m < i < j \leqslant N} Q_{i,j}, \qquad \xi = \sum_{1 \leqslant i \leqslant m} \xi_i,$$

$$T_A = \sum_{1 \leqslant i \leqslant m} T_i, \qquad T_D = \sum_{m < i \leqslant N} T_i,$$

and where the $\xi_i$ are given by (4.6). Furthermore, write $W = T_D + Q_D$. We have $T + Q = T_A + Q_A + \xi + W$ and $\mathrm{e}\{T + Q\} = v\mathrm{e}\{Q_A\}$, with $v = \mathrm{e}\{W + T_A + \xi\}$. Expanding in powers of $iQ_A$ and using symmetry, we obtain

$$\mathrm{Ee}\{T + Q\} = f_1^* + f_2^* + R, \qquad F_1^* = \mathrm{E}v, \qquad f_2^* = \mathrm{i}\binom{m}{2}\mathrm{E}vQ_{1,2}, \qquad (4.14)$$

with $|R| \leqslant \mathrm{E}Q_A^2$. By symmetry, we have

$$\mathrm{E}Q_A^2 = \binom{m}{2}p^2q^2t^2\mathrm{E}y_{1,2}^2 \leqslant \mu^2t^2N^{-1}\gamma_2 \prec \mathscr{R}.$$

Now (4.14) implies $\mathrm{e}\{T + Q\} \sim f_1^* + f_2^*$.

The rest of the proof consists of two steps. In the first step we show that

$$f_2^* \sim h_3, \qquad h_3 = \mathrm{i}^3\binom{m}{2}\mathrm{Ee}\{T^{(1,2)}\}V. \qquad (4.15)$$

In the second step we prove

$$f_1^* \sim h_1 + h_4, \qquad \text{where } h_4 = h_2 - h_3. \qquad (4.16)$$

*Step 1.* We start by showing

$$f_2^* \sim f_3^*, \qquad f_3^* = \mathrm{i}\binom{m}{2}\mathrm{E}v_1Q_{1,2}, \qquad v_1 = \mathrm{e}\{W + T_A + \xi^{(1,2)}\}. \qquad (4.17)$$

Write $v = v_1\mathrm{e}\{\xi_1 + \xi_2\}$. expanding the exponent in powers of $(\xi_1 + \xi_2)$, we obtain

$$f_2^* = f_3^* + f_4^* + f_5^* + f_6^*, \qquad f_j^* = \mathrm{i}^2\binom{m}{2}\mathrm{E}v_1Q_{1,2}l_j, \quad j = 4, 5, 6,$$

$$l_4 = \xi_1 + \xi_2, \qquad l_5 = (\xi_1^2 + \xi_2^2)v_2, \qquad l_6 = 2\xi_1\xi_2v_2, \qquad v_2 = \mathrm{i}^2\mathrm{e}\{\theta_1(\xi_1 + \xi_2)\}(1 - \theta).$$

In order to prove (4.17) we shall show $f_i^* \sim 0$, for $i = 4, 5, 6$.

To show $f_6^* \sim 0$, we bound $|v_1v_2| \leqslant 1$ and obtain

$$|f_6^*| \leqslant m^2\mathrm{E}|Q_{1,2}\xi_1\xi_2| = m^2p^2q^2|t|^3\mathrm{E}|y_{1,2}\zeta_m(A_1)\zeta_m(A_2)|.$$

Combining the inequalities $|\zeta_m(A_1)\zeta_m(A_2)| \leqslant \zeta_m^2(A_1) + \zeta_m^2(A_2)$ and

$$\mathrm{E}|y_{1,2}|\zeta_m^2(A_i) = pq(n - m)\mathrm{E}|y_{1,2}|y_{i,n}^2 \leqslant qN^{-7/2}\gamma_3, \qquad i = 1, 2,$$

and the bound $|t| \leqslant N^{1/2}$, we obtain $|f_6| \ll \mu^2t^2\gamma_3N^{-1} \prec \mathscr{R}$.

Let us show $f_5^* \sim 0$. By symmetry, it suffices to show $m^2\mathrm{E}v_1v_2Q_{1,2}\xi_1^2 \sim 0$. Expanding the exponent in $v_2$ in powers of $\mathrm{i}\theta\xi_2$ and then the exponent in $v_1$ in powers of $\mathrm{i}T_2$ we obtain

$$|\mathrm{E}v_1v_2Q_{1,2}\xi_1^2| \leqslant R_1 + R_2, \qquad R_1 = \mathrm{E}|Q_{1,2}\xi_2|\xi_1^2, \qquad R_2 = \mathrm{E}|Q_{1,2}T_2|\xi_1^2.$$

Invoking (4.45) and the inequality $|t| \leqslant N^{1/2}$, we obtain

$$m^2R_1 \ll \mu^2t^4N^{-5/2}\gamma_4 \prec \mathscr{R}, \qquad m^2R_2 \prec \mu^2t^2N^{-1}(1 + \gamma_4) \prec \mathscr{R},$$

thus completing the proof of $f_5^* \sim 0$.

Let us show $f_4^* \sim 0$. By symmetry, it suffices to show $m^2R \sim 0$ with $R = \mathrm{E}v_1Q_{1,2}\xi_1$. Expanding $v_1$ in powers of $\mathrm{i}T_2$, we can replace $v_1$ by $\mathrm{i}T_2v_3$, with $v_3 = \mathrm{e}\{W + T_A^{(2)} + \xi^{(1,2)} + \theta_1T_2\}$. Now, using the simple bound $|\mathrm{E}_{\{B\}}v_3| \leqslant Z_B$, we obtain

$$|R| \leqslant \mathrm{E}R_1 R_2, \qquad R_1 = |Q_{1,2} T_2|, \qquad R_2 = \mathrm{E}_{[1,2]} |\xi_1| Z_B. \tag{4.18}$$

First we bound $R_2$. By Hölder's inequality,

$$R_2 \leqslant R_3 R_4, \qquad R_3^2 = \mathrm{E}_{[1,2]} \xi_1^2, \qquad R_4^2 = \mathrm{E}_{[1,2]} Z_B^2. \tag{4.19}$$

Furthermore, by (4.9), $R_4^2 \leqslant 2R_5^2 + 2R_6^2$, where

$$R_5^2 = \mathrm{E}_{[1,2]} I_B^2 \leqslant \delta^{-1} \mathrm{E}_{[1,2]} \kappa_B \ll \mathrm{E}_{[1,2]} \kappa_B, \qquad R_6^2 = \mathrm{E}_{[1,2]} \Psi_B^2 \leqslant F_B^2. \tag{4.20}$$

Combining (4.20) with the relations (which follow from symmetry and (4.5))

$$\mathrm{E}_{[1,2]} \xi_1^2 = pq(n-m)(\nu_1 - p)^2 t^2 \mathrm{E}^{1,2} y_{1,n}^2 = q(\nu_1 - p)^2 t^2 N \mathrm{E}^{1,2} y_{1,n}^2, \tag{4.21}$$

$$\mathrm{E}_{[1,2]} \kappa_b = \alpha N \mathrm{E}_{[1,2]} \zeta_m (A_4)^2 = \alpha N(n-m)pq \mathrm{E}^{1,2} y_{4,n}^2 \ll qN^{-1}\gamma_2, \tag{4.22}$$

we obtain

$$R_3 R_5 \ll q|\nu_1 - p||t| \gamma_2^{1/2} \tilde{\gamma}_2^{1/2} \quad \text{and} \quad R_3 R_6 \ll N^{1/2} q^{1/2} |\nu_1 - p||t| \tilde{\gamma}_2^{1/2} F_B. \tag{4.23}$$

Here we denote $\tilde{\gamma}_2 = \mathrm{E}^{1,2} y_{1,n}^2$. using the first inequality of (4.23), we obtain

$$m^2 \mathrm{E}R_1 R_3 R_5 \ll m^2 p^2 q^3 t^2 \gamma_2^{1/2} \mathrm{E}|y_{1,2}| \tilde{\gamma}_2^{1/2} |z_2|,$$

and invoking the second inequality of (4.47), we obtain

$$m^2 \mathrm{E}R_1 R_3 R_5 \ll \mu^2 t^2 N^{-1} \gamma_2^{3/2} \prec \mathscr{R}, \tag{4.24}$$

Using the second inequality of (4.23), we obtain

$$m^2 \mathrm{E}R_1 R_3 R_6 \ll m^2 p^2 q^{5/2} N^{1/2} F_B t^2 \mathrm{E}|y_{1,2}| \tilde{\gamma}_2^{1/2} |z_2|,$$

and invoking the first inequality of (4.47) and (4.10), we obtain

$$m^2 \mathrm{E}R_1 R_3 R_6 \ll \mu^2 N^{-1} F_B t^2 (|t| q^{1/2} + |s|) \gamma_2 \prec \mathscr{R}. \tag{4.25}$$

Since, by (4.18) and (4.19), $|R| \ll \mathrm{E}R_1 R_3 R_5 + \mathrm{E}R_1 R_3 R_6$, it follows from (4.24) and (4.25) that $m^2 R \sim 0$.

In the next step we show that

$$f_3^* \sim f_7^*, \qquad f_7^* = \mathrm{i}^3 \binom{m}{2} \mathrm{E} v_4 V, \qquad v_4 = \mathrm{e}\{w + T_A^{(1,2)} + \xi^{(1,2)}\}. \tag{4.26}$$

Substitute $v_1 = v_4 \mathrm{e}\{T_1 + T_2\}$ in $f_3^*$. Furthermore, using the expansion

$$\mathrm{e}\{T_1 + T_2\} = (1 + T_2 + T_2^2 \mathrm{e}\{\theta_1 T_2\}(1 - \theta_1)) \mathrm{e}\{T_1\}$$

$$= \mathrm{e}\{T_1\} + T_2(1 + T_1 + T_1^2 \mathrm{e}\{\theta_2 T_1\}(1 - \theta_2)) + T_2^2 \mathrm{e}\{\theta_1 T_2\}(1 - \theta_1)(1 + T_1 \mathrm{e}\{\theta_3 T_1\}), \tag{4.27}$$

we obtain $\mathrm{E}v_1 Q_{1,2} = \mathrm{E}v_4 V + R_1 + R_2$, with $|R_i| \leqslant \mathrm{E}Z_B|VT_i|$, $i = 1, 2$. Therefore, in order to prove (4.26) it remains to show $m^2 R_i \sim 0$, for $i = 1, 2$. By symmetry, it suffices to show $m^2 R_1 \sim 0$.

It follows from (4.9) and (4.22) that

$$|R_1| \leqslant \mathrm{E}|VT_1|\kappa_B + \mathrm{E}|VT_1|\Psi_B \ll (N^{-1}q\gamma_2 + F_b)\mathrm{E}|VT_1|. \tag{4.28}$$

Combining (4.46) and the inequalities $|t| \leqslant N^{1/2}$ and $|s| \leqslant (Nq)^{1/2}$, we obtain

$$\mathrm{E}|VT_1| \ll p^2q^2|t|(|t| + |s|q^{-1/2})N^{-2}(\beta_4 + \gamma_4)^{1/2}.$$

Therefore,

$$m^2N^{-1}q\gamma_2\mathrm{E}|VT_1| \ll \mu^2(t^2 + s^2)N^{-1}(\beta_4 + \gamma_4) \prec \mathcal{R}.$$

Finally, (4.46) in combination with (4.10) yields $m^2F_B\mathrm{E}|VT_1| \prec \mathcal{R}$, and this inequality, in view of (4.28), completes the proof of (4.26).

Now we show

$$f_7^* \sim f_8^*, \qquad f_8^* = \mathrm{i}^3\binom{m}{2}\mathrm{E}v_5V, \qquad v_5 = \mathrm{e}\{W + T_A^{(1,2)}\}, \tag{4.29}$$

Expanding $v_4$ in powers of $\mathrm{i}\xi^{(1,2)}$, we obtain

$$\mathrm{E}v_4V = \mathrm{E}v_5V + \mathrm{i}\mathrm{E}v_5V\xi^{(1,2)} + R, \qquad \text{with } |R| \leqslant \mathrm{E}(\xi^{(1,2)})|V|. \tag{4.30}$$

Write $|R| \leqslant \mathrm{E}|V|\mathrm{E}_{[1,2]}(\xi^{(1,2)})^2$. By symmetry and (4.5),

$$\mathrm{E}_{[1,2]}(\xi^{(1,2)})^2 = p^2q^2(m-2)(n-m)t^2\mathrm{E}^{1,2}y_{3,n}^2 \ll \mu qt^2N^{-1}\gamma_2.$$

Now invoking (4.31) – see below – and the bound $|t| \leqslant N^{1/2}$, we obtain

$$m^2|R| \ll \mu^3t^2(t^2 + s^2)N^{-1}\gamma_2(\gamma_2 + 1) \prec \mathcal{R}.$$

This inequality, together with (4.30), implies $f_7^* \sim f_8^* + f_9^*$, where

$$f_9^* = \mathrm{i}^4\binom{m}{2}\mathrm{E}v_5V\xi^{(1,2)} = \mathrm{i}^4\binom{m}{2}(m-2)\mathrm{E}v_5V\xi_3,$$

in symmetry. In order to prove (4.29) it remains to show $f_9^* \sim 0$, to which we now turn.

Expanding $v_5$ in powers of $\mathrm{i}T_3$ we obtain $|\mathrm{E}v_5V\xi_3| \leqslant \mathrm{E}Y_B|V\xi_3T_3|$. Now, using (4.9), we obtain

$$|f_9^*| \leqslant m^3\mathrm{E}|v_5V\xi_3T_3| \leqslant m^3F_B\mathrm{E}|V\xi_3T_3| \leqslant m^3F_B\mathrm{E}|V|\mathrm{E}^{1,2}|\xi_3T_3|.$$

Finally, invoking (4.10) and the bounds (which follow from symmetry and (4.5))

$$\mathrm{E}|V| \leqslant (\mathrm{E}Q_{1,2}^2)^{1/2}(\mathrm{E}T_1^2T_2^2)^{1/2} \leqslant p^2q^2|t|(t^2 + s^2/q)N^{-5/2}\gamma_2^{1/2}, \tag{4.31}$$

$$\mathrm{E}^{1,2}|\xi_3T_3| \leqslant (\mathrm{E}^{1,2}\xi_3^2)^{1/2}(\mathrm{E}^{1,2}T_3^2)^{1/2} \leqslant pq|t|(|t| + |s|)N^{-3/2}\gamma_2^{1/2},$$

we obtain $f_9^* \prec \mathcal{R}$.

Let us show that $f_8^* \sim h_3$. Expanding $v_5$ in powers of $\mathrm{i}Q_D$, we obtain

$$\mathrm{E}v_5V = \mathrm{E}v_6V + \mathrm{i}\mathrm{E}v_6VQ_D + R, \qquad v_6 = \mathrm{e}\{T^{(1,2)}\}, \tag{4.32}$$

with $|R| \leqslant \mathrm{E}Y_B|V|Q_D^2$. Note that, by symmetry,

$$\mathrm{E}Y_B|V|Q_D^2 = \binom{n-m}{2}t^2p^2q^2\mathrm{E}Y_B|V|y_{n-1,n}^2. \tag{4.33}$$

Invoking (4.9) and then using (4.5), we obtain

$$\mathrm{E}Y_B|V|y_{n-1,n}^2 \leqslant F_B\mathrm{E}|V|y_{n-1,n}^2 \leqslant F_B N^{-3}\gamma_2\mathrm{E}|V|. \tag{4.34}$$

Combining (4.33) and (4.34) and then invoking (4.31) and (4.10), we obtain $m^2 R \prec \mathscr{R}$. Now it follows from (4.32) that $f_8^* \sim h_3 + f_{10}^*$, where

$$f_{10}^* = \mathrm{i}^4 \binom{m}{2} \mathrm{E}v_6 VQ_D = \mathrm{i}^4 \binom{m}{2}(n-m)\mathrm{E}v_6 VQ_{n-1,n},$$

by symmetry.

We complete the proof of (4.15) by showing that $f_{10}^* \sim 0$. Expanding $v_6$ in powers of $\mathrm{i}T_{n-1}$ and $\mathrm{i}T_N$, we obtain $|\mathrm{E}v_6 VQ_{n-1,n}| \ll \mathrm{E}|V^*V|Y_B$, where we denote $V^* = T_{n-1}T_NQ_{n-1,n}$. Furthermore, using (4.9) and then invoking the simple inequality $\mathrm{E}^{1,2}|V^*| \ll \mathrm{E}|V^*|$, we obtain

$$\mathrm{E}|V^*V|Y_B \ll F_B\mathrm{E}|V^*V| \ll F_B\mathrm{E}|V^*|\mathrm{E}|V| = F_B(\mathrm{E}|V|)^2.$$

Therefore, $|f_{10}^*| \leqslant m^2(n-m)F_B(\mathrm{E}|V|)^2$. Finally, an application of (4.31) and (4.10) yields $f_{10}^* \prec \mathscr{R}$, thus completing the proof of (4.15).

*Step 2.* In order to prove (4.16) it suffices to show that

$$f_1^* \sim f_{11}^* + f_{12}^*, \qquad f_{11}^* = \mathrm{E}v_7, \qquad f_{12}^* = \mathrm{i}\mathrm{E}v_7\xi, \qquad v_7 = \mathrm{e}\{T + Q_D\}, \quad (4.35)$$

$$f_{11}^* \sim h_1 + f_{13}^*, \qquad f_{13}^* = \binom{n-m}{2}\mathrm{i}^3\mathrm{E}\mathrm{e}\{T^{(1,2)}\}V, \tag{4.36}$$

$$f_{12}^* \sim f_{14}^*, \qquad f_{14}^* = m(n-M)\mathrm{i}^3\mathrm{E}\mathrm{e}\{t^{(1,2)}\}V. \tag{4.37}$$

We begin by expanding $v$ in powers of $\mathrm{i}\xi$, to obtain

$$f_1^* = f_{11}^* + f_{12}^* + f_{15}^*, \qquad \text{with } f_{15}^* = \mathrm{i}^2\mathrm{E}v_7\xi^2\mathrm{e}\{\theta_1\xi\}(1-\theta_1).$$

In order to prove (4.35), we shall show $f_{15}^* \sim 0$. Split

$$\Omega_m = S_1 \cup S_2 \cup S_3 \cup S_4,$$

with $S_i \cap S_j = \varnothing$, $i \neq j$, and $|S_j| \approx m/4$, $1 \leqslant j \leqslant 4$. Split $\xi = \delta_1 + \ldots + \delta_4$, where $\delta_j = \sum_{i \in S_j}\xi_i$. We have

$$f_{15}^* = \sum_{1 \leqslant j,k \leqslant 4} r_{j,k}, \qquad r_{j,k} = \mathrm{i}^2\mathrm{E}v_7\delta_j\delta_k\mathrm{e}\{\theta_1\xi\}(1-\theta_1).$$

We shall show $r_{j,k} \sim 0$, for every $1 \leqslant j,\ k \leqslant 4$. By symmetry, it suffices to prove $r_{1,1} \sim 0$ and $r_{1,2} \sim 0$.

Let us show $r_{1,1} \sim 0$. expanding in powers of $\mathrm{i}\theta_1\delta_2$, we obtain

$$\mathrm{e}\{\theta_1\xi\} = v_8 + \mathrm{i}\delta_2 v_8\tilde{v}, \qquad v_8 = \mathrm{e}\{\theta_1(\delta_1 + \delta_3 + \delta_4)\}, \qquad \tilde{v} = \theta_1\mathrm{E}_{\theta_2}\mathrm{e}\{\theta_1\theta_2\delta_2\}.$$

Substitution of this formula gives

$$r_{1,1} = R_1 + R_2, \qquad R_1 = \mathrm{i}^2\mathrm{E}v_7v_8\delta_1^2(1-\theta_1), \qquad R_2 = \mathrm{i}^3\mathrm{E}v_7v_8\tilde{v}\delta_1^2\delta_2(1-\theta_1).$$

Similarly, expanding $v_8$ in powers of $i\theta_1\delta_3$, we obtain $R_2 = R_3 + R_4$, where

$$R_3 = i^3 E v_7 v_9 \delta_1^2 \delta_2 \tilde{v}(1 - \theta_1), \qquad v_9 = e\{\theta_1(\delta_1 + \delta_4)\}, \qquad |R_4| \leqslant E\delta_1^2|\delta_2\delta_3|.$$

Therefore, $|r_{1,1}| \leqslant |R_1| + |R_2| + |R_3|$. Furthermore, invoking the inequalities $|E_{\{s_2\}}v_7v_8| \leqslant Y_{S_2}$ and $|E_{\{S_3\}}v_7v_9\tilde{v}| \leqslant Y_{S_3}$, we obtain

$$|r_{1,1}| \leqslant r_1 + r_2 + r_3, \qquad r_1 = E\delta_1^2 Y_{S_2}, \qquad r_2 = E\delta_1^2|\delta_2|Y_{S_3}, \qquad r_3 = E\delta_1^2|\delta_2\delta_3|. \quad (4.38)$$

Now we show $r_i \sim 0$, for $i = 1, 2, 3$. Denote for brevity $m_i = |S_i|$, $1 \leqslant |i| \leqslant 4$.

Let us show $r_2 \sim 0$. By symmetry,

$$E_{\{S_1\}}\delta_1^2 = m_1 pqt^2\zeta_m^2(A_1), \qquad E_{\{S_2\}}\delta_2^2 = m_2 pqt^2\zeta_m^2(A_{i_0}), \quad (4.39)$$

with $i_0 \in S_2$. Combining (4.39) and the inequality $E_{\{S_2\}}|\delta_2| \leqslant (E_{\{S_2\}}\delta_2^2)^{1/2}$ and using symmetry again, we obtain

$$r_2 = EY_{S_3}(E_{\{S_1\}}\delta_1^2)E_{\{S_2\}}|\delta_2| \leqslant m_1 m_2^{1/2}(pq)^{3/2}|t|^3 E\zeta_m^2(A_1)|\zeta_m(A_{i_0})|Y_{S_3}$$

$$\ll m^{3/2}(pq)^{3/2}|t|^3 E|\zeta_m(A_1)|^3 Y_{S_3}. \quad (4.40)$$

In the last step we applied (4.44) and again used symmetry. Furthermore, invoking (5.4) and using symmetry and (4.9), we obtain

$$E|\zeta_m(A_1)|^3 Y_{S_3} \ll N^{1/2} pq(n - m)E|y_{1,n}|^3 Y_{S_3} \leqslant N^{-3} F_{S_3}\gamma_3.$$

this inequality, in combination with (4.40) and (4.10), implies $r_2 \sim 0$.

To show $r_1 \sim 0$ we use symmetry, and apply (4.9) and (4.10):

$$r_1 = m_1 t^2 p^2 q^2(n - m)Ey_{1,n}^2 Y_{S_2} \ll t^2 F_{S_2} N^{-1}\gamma_2 \prec \mathscr{R}. \quad (4.41)$$

To show $r_3 \sim 0$, we first use (4.44) to obtain $r_3 \leqslant E\delta_1^2\delta_2^2 + E\delta_1^2\delta_3^2$ and then apply (4.48). Finally, collecting the bounds $r_i \prec \mathscr{R}$, $i = 1, 2, 3$ in (4.38) we get $r_{1,1} \sim 0$.

Let us show $r_{1,2} \sim 0$. Expanding in powers of $i\theta_1\delta_3$ and $i\theta_1\delta_4$, we obtain

$$e\{\theta_1\xi\} = v_{10} + v_{10}v_{11}i\theta_1\delta_3, \qquad v_{10} = e\{\theta_1(\delta_1 + \delta_2 + \delta_4)\}, \qquad v_{11} = E_{\theta_2}e\{\theta_1\theta_2\delta_3\},$$

$$v_{10} = v_{12} + v_{12}v_{13}i\theta_1\delta_4, \qquad v_{12} = e\{\theta_1(\delta_1 + \delta_2)\}, \qquad v_{13} = E_{\theta_3}e\{\theta_1\theta_3\delta_4\}.$$

Combining these expansions, we obtain

$$e\{\theta_1\xi\} = v_{10} + v_{11}v_{12}i\theta_1\delta_3 + v_{11}v_{12}v_{13}i^2\theta_1^2\delta_3\delta_4.$$

The last identity, in combination with the bounds $|E_{\{S_3\}}v_7v_{10}| \leqslant Y_{S_3}$ and $|E_{\{S_4\}}v_7v_{11}v_{12}| \leqslant Y_{S_4}$, implies

$$r_{1,2} \leqslant E|\delta_1\delta_2|Y_{S_3} + E|\delta_1\delta_2\delta_3|Y_{S_4} + E|\delta_1\delta_2\delta_3\delta_4|$$

$$\leqslant E\delta_1^2 Y_{S_3} + E\delta_2^2 Y_{S_3} + E\delta_1^2|\delta_2|Y_{S_4} + E\delta_3^2|\delta_2|Y_{S_4} + E\delta_1^2\delta_2^2 + E\delta_3^2\delta_4^2. \quad (4.42)$$

In the last step we used the simple inequality $ab \leqslant a^2 + b^2$ several times. Note that the quantities in (4.42) can be bounded in the same way as $r_1$, $r_2$, and $r_3$ above in the proof of $r_{1,1} \sim 0$. Hence, $r_{1,2} \sim 0$ and this completes the proof of (4.35).

Let us prove (4.36). Expanding $v_7$ in powers of $iQ_D$, we obtain

$$f_{11}^* = h_1 + f_{16}^* + R, \qquad f_{16}^* = iEe\{T\}Q_D,$$

with $|R| \leqslant EY_B Q_D^2$. Furthermore, by symmetry,

$$f_{16}^* = \binom{n-m}{2} iEe\{T\}Q_{1,2} \quad \text{and} \quad EY_B Q_D^2 = \binom{n-m}{2} p^2 q^2 t^2 EY_B y_{n-1,n}^2.$$

Combining (4.9) and (4.10), we obtain $R \prec \mathscr{R}$ and, therefore, $f_{11}^* \sim h_1 + f_{16}^*$.

Let us show $f_{16}^* \sim f_{13}^*$. Write $e\{T\} = e\{T^{(1,2)}\}e\{T_1 + T_2\}$ and use (4.27) to obtain

$$Ee\{T\}Q_{1,2} = i^2 Ee\{T^{(1,2)}\}V + R_1 + R_2, \qquad \text{with } |R_i| \ll E|VT_i|Y_B. \qquad (4.43)$$

By (4.9), $|R_i| \ll F_B E|VT_i|$. Furthermore, invoking (4.46) and (4.10) we obtain $n^2 R_i \prec \mathscr{R}$, $i = 1, 2$. These bounds, together with (4.43), imply $f_{16}^* \sim f_{13}^*$, thus completing the proof of (4.36).

Let us prove (4.37). By symmetry, $f_{12}^* = miEv_7\xi_1$. Expanding $v_7$ in powers of $iT_1$, we obtain

$$f_{12}^* = f_{17}^* + R_1, \qquad f_{17}^* = mi^2 Ee\{T^{(1)} + Q_D\}\xi_1 T_1, \qquad |R_2| \leqslant mEY_B|\xi_1|T_1^2.$$

Furthermore, expanding the exponent in powers of $iQ_D$, we obtain

$$f_{17}^* = f_{18}^* + R_2, \qquad f_{18}^* = mi^2 Ee\{T^{(1)}\}\xi_1 T_1, \qquad |R_2| \leqslant mEY_B|\xi_1 T_1 Q_D|.$$

Note, that by symmetry, $f_{18}^* = m(n-m)i^2 Ee\{T^{(1)}\}Q_{1,2}T_1$. Finally, expanding the exponent in powers of $iT_2$, we obtain

$$f_{18}^* = f_{14}^* + R_3, \qquad \text{with } |R_3| \leqslant n(n-m)EY_B|VT_2|.$$

Therefore in order to prove (4.37) it remains to show $R_i \prec \mathscr{R}$, for $i = 1, 2, 3$.

To show $R_1 \prec \mathscr{R}$ use the inequality $|\xi_1|T_1^2 \leqslant \xi_1^2 + T_1^4$. We obtain $|R_1| \leqslant R_{1,1} + R_{1,2}$, with $R_{1,1} = mEY_B T_1^4$ and $R_{1,2} = mEY_B\xi_1^2$. By (4.9) and (4.10), $R_{1,1} \prec \mathscr{R}$. Furthermore, the bound $R_{1,2} \prec \mathscr{R}$ is obtained in the same way as (4.41).

To show $R_2 \prec \mathscr{R}$ use the inequality $|\xi_1 T_1 Q_D| \leqslant \xi_1^2 + T_1^2 Q_D^2$. We get $|R_2| \leqslant R_{2,1} + R_{2,2}$, with $R_{2,1} = mEY_B\xi_1^2 \prec \mathscr{R}$ (cf. (4.41)) and with

$$R_{2,2} = mEY_B T_1^2 Q_D^2 = m\binom{n-m}{2}EY_B T_1^2 Q_{n-1,n}^2 \leqslant mn^2 F_B E T_1^2 Q_{n-1,n}^2,$$

by symmetry and (4.9). Now, combining (4.10) and the inequality

$$ET_1^2 Q_{n-1,n}^2 = p^3 q^3 t^2 Ez_1^2 E^1 y_{n-1,n}^2 \ll p^3 q^3 t^2(t^2 + s^2/q)N^{-4}\gamma_2,$$

(here we use (4.5)) we obtain $R_{2,2} \prec \mathscr{R}$.

To show $R_3 \prec \mathscr{R}$ we apply (4.9) to obtain $R_3 \leqslant nmF_B E|VT_2|$. Then combining (4.46) and (4.10) we obtain $R_3 \prec \mathscr{R}$. We arrive at (4.37), thus completing the proof of the lemma. $\qquad \square$

In the next lemma we gather together some auxiliary inequalities used in Lemma 4.1. We shall frequently use the inequalities

$$ab \leqslant a^2 + b^2, \qquad a^2 b \leqslant a^3 + b^3. \tag{4.44}$$

**Lemma 4.2** *We have*

$$\mathrm{E}|Q_{1,2}\xi_2|\xi_1^2 \ll p^2 q^3 t^4 N^{-9/2}\gamma_4, \qquad \mathrm{E}|Q_{1,2}T_2|\xi_1^2 \ll p^2 q^3 |t|^3 N^{-7/2}(1 + \gamma_4), \tag{4.45}$$

$$\mathrm{E}|Q_{1,2}T_2 T_1^2| \ll p^2 q^2 |t|(|t|^3 + |s|^3 q^{-3/2})N^{-3}(\beta_4 + \gamma_4)^{1/2}, \tag{4.46}$$

$$\mathrm{E}|y_{1,2}x_2|\tilde{\gamma}_2^{1/2} \ll N^{-7/2}\gamma_2, \qquad \mathrm{E}|y_{1,2}z_2|\tilde{\gamma}_2^{1/2} \ll N^{-3}\gamma_2, \qquad \tilde{\gamma}_2 = \mathrm{E}^{1,2}y_{1,n}^2, \tag{4.47}$$

$$\mathrm{E}\delta_K^2 \delta_M^2 \ll p^2 q^2 t^2 m^2 N^{-3}\gamma_4, \qquad \text{for any } K, M \subset \Omega_m, K \cap M = \varnothing, \tag{4.48}$$

*with* $|K|, |M| > 0$. *Here* $\delta_K = \sum_{i \in K} \xi_i$.

***Proof.*** Let us prove (4.45). We have

$$\mathrm{E}|Q_{1,2}\xi_2\delta_1^2| \ll p^2 q^2 t^4 \mathrm{E}|y_{1,2}\zeta_m(A_2)|\zeta_m^2(A_1) \leqslant 2p^2 q^2 t^4 \mathrm{E}|y_{1,2}\zeta_m^3(A_1)|,$$

where in the last step we apply (4.44) and use symmetry. Furthermore, writing $\mathrm{E}|y_{1,2}\zeta_m^3(A_1)| = \mathrm{E}|y_{1,2}|\mathrm{E}^{1,2}|\zeta_m^3(A_1)|$ and invoking (5.4), we obtain the first inequality in (4.45). To prove the second inequality we apply (4.21),

$$\mathrm{E}|Q_{1,2}T_2|\xi_1^2 = \mathrm{E}|Q_{1,2}T_2|\mathrm{E}_{[1,2]}\xi_1^2 \leqslant p^2 q^3 N|t|^3 \mathrm{E}y_{1,n}^2|y_{1,2}z_2|,$$

and use the inequality $\mathrm{E}y_{1,n}^2|y_{1,2}z_2| \leqslant N^{-9/2}(1 + \gamma_4)$. To prove this inequality combine (4.5), Hölder's inequality and the bounds

$$|y_{1,2}z_2| \leqslant \sqrt{N}|y_{1,2}x_2| + |y_{1,2}|, \qquad |y_{1,2}x_2| \leqslant N y_{1,2}^2 + N^{-1}x_2^2.$$

Let us prove (4.46). We have $\mathrm{E}|Q_{1,2}T_2|T_1^2 \leqslant p^2 q^2 |t|\mathrm{E}|y_{1,2}z_2|z_1^2$. Now (4.46) follows from the bound

$$\mathrm{E}|y_{1,2}z_2|z_1^2 \leqslant (|t|^3 + t^2|s|q^{-1/2} + |t|s^2/q + |s|^3 q^{-3/2})N^{-3}(\beta_4 + \gamma_4)^{1/2}$$

which is a consequence of the following inequalities:

$$z_1^2|z_2| \leqslant (t^2 x_1^2 + s^2/\tau^2)|z_2|, \qquad \mathrm{E}|y_{1,2}x_2| \leqslant (\mathrm{E}y_{1,2}^2)^{1/2}(\mathrm{E}x_2^2)^{1/2} \leqslant N^{-2}\gamma_2^{1/2},$$

$$\mathrm{E}|y_{1,2}x_2|x_1^2 \leqslant (\mathrm{E}y_{1,2}^2 x_1^2)^{1/2}(\mathrm{E}x_1^2 x_2^2)^{1/2} \leqslant N^{-1}(\mathrm{E}y_{1,2}^4 \mathrm{E}x_1^4)^{1/4} \leqslant N^{-3}(\beta_4 + \gamma_4)^{1/2},$$

$$\mathrm{E}|y_{1,2}|x_1^2 \leqslant (\mathrm{E}y_{1,2}^2 x_1^2)^{1/2}(\mathrm{E}x_1^2)^{1/2} \leqslant N^{-1/2}(\mathrm{E}y_{1,2}^4 \mathrm{E}x_1^4)^{1/2} \leqslant N^{-5/2}(\beta_4 + \gamma_4)^{1/2}.$$

Let us prove (4.47). By (4.44), $|y_{1,2}x_2|\tilde{\gamma}_2^{1/2} \leqslant N^{-1/2}y_{1,2}^2 + N^{1/2}x_2^2\tilde{\gamma}_2$. Now invoking (4.5), we obtain the first inequality of (4.47). To prove the second one, write

$$\mathrm{E}|y_{1,2}z_1|\tilde{\gamma}_2^{1/2} \leqslant N^{1/2}\mathrm{E}|y_{1,2}x_2|\tilde{\gamma}_2^{1/2} + \mathrm{E}|y_{1,2}|\tilde{\gamma}_2^{1/2}$$

(where we have used the bounds $|t| \leqslant N^{1/2}$ and $|s| \leqslant \tau$) and apply Hölder's inequality to the second summand.

Let us prove (4.48). By symmetry,

$$\mathrm{E}\delta_K^2\delta_M^2 = |K\|M|p^2q^2t^4\mathrm{E}\zeta_m^2(A_1)\zeta_m^2(A_2). \tag{4.49}$$

A simple calculation yields

$$\mathrm{E}\zeta_m^2(A_1)\zeta_m^2(A_2) \leqslant pq(n-m)N^{-6}\gamma_4 + p^2q^2(n-m)(n-m-1)N^{-6}\gamma_2^2 \ll N^{-4}\gamma_4.$$

Substituting this bound in (4.49) and using the inequalities $|K|, |M| < m$ and $t^2 \leqslant N$, we obtain (4.48). $\qquad\square$

# 5. Auxiliary results

**Lemma 5.1.** *For the random variables $v_i$ and $\Theta_i$, $i = 1, 2, 3$, defined in (3.6) and (3.25) above, the following inequalities hold:*

$$\mathrm{E}\Theta_i^2 \leqslant N^{-2}\gamma_2, \qquad i = 1, 2, 3, \tag{5.1}$$

$$\mathrm{E}^*|v_i(A_0^*)| \leqslant q^{-1}\Theta_i, \qquad i = 1, 2. \tag{5.2}$$

*For $\Lambda_m = \sum_{1 \leqslant i < j \leqslant m} g_2(A_i, A_j)$, with $3 \leqslant m \leqslant n$, we have*

$$\mathrm{E}\Lambda_m^2 = \frac{m(m-1)}{2N^3}(1 - c_\Lambda)\gamma_2, \qquad c_\Lambda = \frac{2(m-2)}{n-2} - \frac{(m-2)(m-3)}{(n-2)(n-3)}. \tag{5.3}$$

*For the random variable $\zeta_m(A_k)$ defined in (4.6), the inequality*

$$\mathrm{E}_{\{m+1,\dots,n\}}|\zeta_m(A_k)|^3 \ll pq(npq)^{1/2} \sum_{j=m+1}^{n} |g_2(A_k, A_j)|^3, \qquad K \leqslant m \leqslant n, \tag{5.4}$$

*holds; recall the definition of $\mathrm{E}_{\{m+1,\dots,n\}}$ just before (4.6).*

**Proof.** We shall prove (5.1) for case $i = 1$ only. For $i = 2, 3$ the proof is similar. By Hölder's inequality,

$$\Theta_1^2 \leqslant \mathrm{E}^*v_1^2(A_1^*) = \sum_{k+1 \leqslant i,j \leqslant N} \mathrm{E}^* g_2(A_1^*, A_i)g_2(A_1^*, A_j).$$

By symmetry,

$$\mathrm{E}\mathrm{E}^* g_2(A_1^*, A_i)g_2(A_1^*, A_j) = \mathrm{E} g_2(A_1, A_i)g_2(A_1, A_j),$$

and therefore,

$$\mathrm{E}\Theta_1^2 \leqslant (N-k)\mathrm{E}g_2^2(A_1, A_2) + (N-k)(N-k-1)\mathrm{E}g_2(A_1, A_2)g_2(A_1, A_3).$$

Now, invoking the identity

$$\mathrm{E}g_2(A_1, A_2)g_2(A_1, A_3) = -(n-2)^{-1}\mathrm{E}g_2^2(A_1, A_2), \tag{5.5}$$

(use (1.2)) we complete the proof of $\mathrm{E}\Theta_1^2 \leqslant N^{-2}\gamma_2$.

To prove (5.2) we combine the obvious inequality $|\mathscr{J}_i|/|\mathscr{J}_0| \leqslant q^{-1}$ and the inequalities

$$\mathrm{E}^*|v_i(A_0^*)| = |\mathcal{J}_0|^{-1} \sum_{k \in \mathcal{J}_0} |v_i(A_k)| \leqslant |\mathcal{J}_0|^{-1} \sum_{k \in \mathcal{J}_i} |v_i(A_k)| \leqslant \Theta_i \frac{|\mathcal{J}_i|}{|\mathcal{J}_0|}, \qquad i = 1, 2.$$

Let us prove (5.3). By symmetry, $\mathrm{E}\Lambda_m^2 = 2^{-1}(m-1)m\mathrm{E}g_2(A_1, A_2)\Lambda_m$. Furthermore,

$$\mathrm{E}g_2(A_1, A_2)\Lambda_m = \mathrm{E}g_2^2(A_1, A_2) + 2(m-2)\mathrm{E}g_2(A_1, A_2)g_2(A_1, A_3)$$

$$+ 2^{-1}(m-2)(m-3)\mathrm{E}g_2(A_1, A_2)g_2(A_3, A_4).$$

Now, invoking (5.5) and the identity

$$\mathrm{E}g_2(A_1, A_2)g_2(A_3, A_4) = 2(n-2)^{-1}(n-3)^{-1}\mathrm{E}g_2^2(A_1, A_2),$$

(use (1.2)) we obtain (5.3).

In order to prove (5.4) we apply Rosenthal's inequality,

$$\mathrm{E}|Z_1 + \ldots + Z_j|^r \leqslant c(r)\sum_{l=1}^j \mathrm{E}|Z_l|^r + c(r)\left(\sum_{l=1}^j \mathrm{E}Z_l^2\right)^{r/2}, \qquad r \geqslant 2,$$

where $Z_1 \ldots, Z_j$ are independent and centred ransom variables. We apply this inequality to the sum $\zeta_m(A_k)$ – cf. (4.15) – conditionally given $\bar{A}$,

$$\mathrm{E}_{\{m+1,\ldots,n\}}|\zeta_m(A_k)|^3 \ll pq \sum_{l=m+1}^n |g_2(A_k, A_l)|^3 + \left(pq \sum_{l=m+1}^n g_2^2(A_k, A_l)\right)^{3/2}.$$

Finally, using Hölder's inequality, we bound the second sum above by

$$\left(\sum_{l=m+1}^n g_2^2(a, A_l)\right)^{3/2} \leqslant (n-m)^{1/2} \sum_{l=m+1}^n |g_2(a, A_l)|^3, \qquad a \in \mathcal{A},$$

thus arriving at (5.4). $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Lemma 5.2.** *For each $0 < d < \pi$ and $x, y \in \mathbb{R}$, and $\beta(x)$ defined in Section 3.1, we have*

$$|\beta(x+y)|^2 \leqslant u_{[d]}(x)v_{[d]}(y), \qquad \text{where } v_{[d]}(y) = 1 + pq\frac{2\pi}{d}\left(\frac{4}{\Theta(d)} + 1\right)y^2,$$

*and where the function $u_{[d]}$ id defined in (4.8).*

**Proof.** In the case where $|x| \geqslant \pi + d$, we have $u_{[d]}(x) = 1$ and the desired inequality follows from the simple bound $|\beta(x+y)| \leqslant 1$.

In the case where $|x| < \pi + d$, we apply the mean value theorem to obtain

$$|\cos(x+y) - \cos(x)| \leqslant |\mathrm{E}\sin(x+\theta_1 y)y|$$

$$\leqslant (|x| + |y|)|y| \leqslant cx^2 + (c^{-1}+1)y^2. \qquad\qquad (5.6)$$

In the last step we applied the inequality $|xy| \leqslant cx^2 + c^{-1}y^2$, with $c > 0$. Combining (5.6) and the indentity $|\beta(x+y)|^2 = 1 - 2pq(1 - \cos(x+y))$, we obtain

$$|\beta(x+y)|^2 \leqslant 1 - 2pq(1 - \cos(x) - cx^2 - (c^{-1}+1)y^2).$$

Now invoking (5.15), we obtain

$$|\beta(x+y)|^2 \leqslant w_1 + w_2, \qquad w_1 = 1 - pq(\Theta(d) - 2c)x^2, \qquad w_2 = (c^{-1}+1)2pqy^2.$$

But $1 - pqx^2\Theta(d) \geqslant d/\pi$, for $|x| \leqslant \pi + d$. Hence, $w_1 > d/\pi$ and therefore $w_1 + w_2 \leqslant w_1(1 + \pi d^{-1}w_2)$. Choosing $c = \Theta(d)/4$ completes the proof of the lemma.  □

**Lemma 5.3.** *Assume that $\beta_2 = 1$. For every $s$, $t \in \mathbb{R}$ and $0 < d < \pi$, we have*

$$\mathbb{E}Z^2(A_1)I_{[d]}(A_1) \geqslant \left(\frac{t^2}{N} + s^2\right)(1 - 2c_{[d]}), \qquad c_{[d]} = \max\left\{\frac{b_1}{d}; \frac{b_1^2}{d^2}\right\}.$$

*Here $Z(a) = tg_1(a) + s$ and $I_{[d]}(a) = I\{H_1|g_1(a)| < d\}$, for $a \in \mathcal{A}$.*

**Remark.** A similar inequality was used by Höglund (1978), where the constatnt (corresponding to $c_{[d]}$) was not specified. For our purposes the dependence of $c_d$ on the parameters $b_1$ and $d$ is important and thus we include the proof.

**Proof.** Denote $\mathcal{K}_{[d]} = \{k : I_{[d]}(a_k) = 0\}$. Clearly, for $r > 0$,

$$|\mathcal{K}_{[d]}| = \sum_{k \in \mathcal{K}_{[d]}} 1 \leqslant \sum_{k \in \mathcal{K}_{[d]}} |g_1(a_k)H_1/d|^r \leqslant n\beta_r\beta_3^{-r}b_1^r d^{-r}. \qquad (5.7)$$

Furthermore, since $\mathbb{E}Z^2(A_1) = t^2 N^{-1} + s^2$, we have

$$\mathbb{E}Z^2(A_1)^2 I_{[d]}(A_1) = \frac{t^2}{N} + s^2 - Wn^{-1}, \qquad W = \sum_{k \in \mathcal{K}_{[d]}} Z^2(a_k). \qquad (5.8)$$

The inequality $(a+b)^2 \leqslant 2a^2 + 2b^2$ implies $W \leqslant 2W_1 + 2W_2$, where

$$W_1 = s^2|K_{[d]}|, \qquad W_2 = t^2 \sum_{k \in \mathcal{K}_{[d]}} g_1^2(a_k) \leqslant \frac{t^2}{N} n^{2/3}\beta_3^{2/3}|\mathcal{K}_{[d]}|^{1/3}.$$

In the last step we applied Hölder's inequality to obtain

$$\sum_{k \in \mathcal{K}_{[d]}} g_1^2(a_k) \leqslant \left(\sum_{k \in \mathcal{K}_{[d]}} |g_1^3(a_k)|\right)^{2/3} |\mathcal{K}_{[d]}|^{1/3}.$$

Now, (5.7) (with $r = 2$) implies $W_1 \leqslant s^2 nc_d$. Furthermore, (5.7) (with $r = 3$) implies $W_2 \leqslant t^2 N^{-1}nc_d$. These inequalities, combined with (5.8), complete the proof.  □

**Lemma 5.4.** *Assume that $\beta_2 = 1$ and that (3.3) holds. For $|t| \leqslant H_1$ and $|s| \leqslant \pi\tau$, the inequalities (4.9) hold true.*

**Proof.** Throughout this proof we use the notation introduced in Section 4. Fix $B \subset \Omega_m$. By Lemma 5.2,

$$Z_B = \prod_{k \in B} |\beta(z_k + t\zeta_m(A_k))| \leqslant \eta_1 \eta_2, \qquad \eta_1^2 = \prod_{k \in B} u_{[1]}(z_k), \qquad \eta_2^2 = \prod_{k \in B} v_{[1]}(t\zeta_m(A_k)).$$

Using the inequality $1 + x \leqslant \exp\{x\}$, we obtain $\eta_2^2 \leqslant \exp\{pqt^2 \kappa_B |B|/N\}$ and therefore, $(1 - I_B)\eta_2 \leqslant g_B(t)$. Finally, combining the inequalities $Z_B \leqslant 1$ and $Z_B \leqslant \eta_1 \eta_2$, we obtain

$$Z_B = I_B Z_B + (1 - I_B) Z_B \leqslant I_B + (1 - I_B)\eta_1 \eta_2 \leqslant I_B + \Psi_B,$$

thus proving the first inequality in (4.9). Here the random variables $\Psi_B = \eta_1 g_B(t)$, $\kappa_B$ and $I_B$ are defined in (4.7) and the function $g_B(t)$ is given by (4.8).

Clearly, Lemma 5.2 (with $x = z_k$ and $y = 0$) implies the second inequality of (4.9).

To prove the third and last one, observe that by Hölder's inequality we have $E^{i_1,\dots,i_4} \Psi_B \ll (E^{i_1,\dots,i_4} \Psi_B^2)^{1/2}$ and thus, it suffices to show

$$E^{i_1,\dots,i_4} \Psi_B^2 \leqslant F_B^2, \qquad \text{for every } i_1, \dots, i_4 \in \Omega_n \backslash B. \tag{5.9}$$

To prove (5.9) note that the inequalities $|t| \leqslant H_1$ and $|s| \leqslant \pi\tau$ imply

$$u_{[1]}(z_k) \leqslant w(A_k), \qquad w(A_k) = 1 - \frac{pq}{2}\Theta(1)z_k^2 I_{[1]}(A_k), \qquad k \in \Omega_n,$$

where we denote $I_{[1]}(a) = I\{H_1 |g_1(a)| < 1\}$. Therefore,

$$\Psi_B^2 = g_B^2(t)\eta_1^2 \leqslant g_B^2(t)\eta, \qquad \text{where } \eta = \prod_{k \in B} w(A_k). \tag{5.10}$$

Denote $D_1 = \{i_1, \dots, i_4\}$ and $D_2 = \Omega_n \backslash D_1$. By Theorem 4 of Hoeffding (1963),

$$E^{i_1,\dots,i_4} \eta \leqslant w_*^{|B|}, \qquad \text{where } w_* = 1 - \frac{pq}{2}\Theta(1)\Gamma_*, \quad \Gamma_* = \frac{1}{|D_2|}\sum_{k \in D_2} z_k^2 I_{[1]}(A_k). \tag{5.11}$$

Below we construct the following lower bound for $\Gamma_*$,

$$\Gamma_* \geqslant \frac{9}{10}\left(t^2 + \frac{s^2}{q}\right)\frac{1}{N}. \tag{5.12}$$

Combining (5.11), (5.12) and the inequality $1 + x \leqslant \exp\{x\}$ we obtain $\eta \leqslant \exp\{-0.45pq\Theta(1)(t^2 + s^2/q)|B|N^{-1}\}$. Now (5.9) follows from (5.10).

Let us prove (5.12). We have

$$\Gamma_* = \frac{n}{n-4}Ez_1^2 I_{[1]}(A_1) - \frac{1}{n-4}M, \qquad M = \sum_{k \in D_1} z_k^2 I_{[1]}(A_k). \tag{5.13}$$

The simple inequality $(a + b)^2 \leqslant 2a^2 + 2b^2$ gives

$$M \leqslant 8s^2/\tau^2 + 2t^2 M_1, \qquad M_1 = \sum_{k \in D_1} g_1^2(A_k). \tag{5.14}$$

By Hölder's inequality and (3.3),

$$M_1 \leqslant 4^{1/3} \left( \sum_{k \in D_1} |g_1(A_k)|^3 \right)^{2/3} \leqslant 4^{1/3} \beta_3^{2/3} n^{2/3} N^{-1} \leqslant (4c_0)^{1/3} \frac{n}{N}.$$

Here we have estimated $\beta_3^2/n \leqslant \beta_4/n \leqslant c_0$; see (3.3). This inequality, in combination with (5.14) implies $M \leqslant 100^{-1} n(t^2 + s^2/q) N^{-1}$, provided that $c_0$ in (3.3) is sufficiently small. Substituting this bound in (5.13) and invoking Lemma 5.3, we obtain (5.12) thus completing the proof of the lemma. ☐

We have used, but not so far stated Hoeffding's (1963) Theorem 4. Consider a population $\mathscr{P}$ of $n$ numbers $p_1, \ldots, p_n$. Let $\mathscr{X}_1, \ldots, \mathscr{X}_N$ denote a random sample *without* replacement from $\mathscr{P}$ and let $\mathscr{Y}_1, \ldots, \mathscr{Y}_N$ denote a random sample *with* replacement from $\mathscr{P}$. In particular, $\mathscr{Y}_1, \ldots, \mathscr{Y}_N$ are independent random variables.

**Theorem (Hoeffding 1963).** *If the function* $f(x)$ *is continuous and convex then*

$$\mathrm{E}f\left( \sum_{k=1}^{N} \mathscr{X}_k \right) \leqslant \mathrm{E}f\left( \sum_{k=1}^{N} \mathscr{Y}_k \right).$$

We conclude this paper by stating two inequalities proved by Höglund (1978):

$$1 - \cos v \geqslant \frac{1}{2} v^2 \Theta(u), \qquad \text{for } |v| \leqslant \pi + u \quad \text{and} \quad 0 \leqslant u \leqslant \pi, \tag{5.15}$$

$$\frac{\pi^{1/2}}{2} \leqslant \binom{n}{N} s^N (1-s)^{n-N} (2\pi s(1-s)n)^{1/2} \leqslant 1, \qquad \text{with } s = \frac{N}{n}, \tag{5.16}$$

where $1 \leqslant N \leqslant n$.

# Acknowledgements

# References

Albers, W., Bickel, P.J. and van Zwet, W.R. (1976) Asymptotic expansions for the power of distribution free tests in the one-sample problem. *Ann. Statist.*, **4**, 108–156.

Babu, G.J. and Singh, K. (1985) Edgeworth expansions for sampling without replacement from finite populations. *J. Multivariate Anal.*, **17**, 261–278.

Bentkus, V., Götze, F. and van Zwet, W.R. (1997) An Edgeworth expansion for symmetric statistics. *Ann. Statist.*, **25**, 851–896.

Bickel, P.J. (1974) Edgeworth expansions in nonparametric statistics. *Ann. Statist.*, **2**, 1–20.

Bickel, P.J. and Robinson, J. (1982) Edgeworth expansions and smoothness. *Ann. Probab.*, **10**, 500–503.

Bickel, P.J. and van Zwet, W.R. (1978) Asymptotic expansions for the power of distribution free tests in the two-sample problem. *Ann. Statist.*, **6**, 937–1004.

Bickel, P.J., Götze, F. and van Zwet, W.R. (1986) The Edgeworth expansion of *U*-statistics of degree two. *Ann. Statist.*, **14**, 1463–1484.

Bikelis, A. (1969) On the estimation of the remainder term in the central limit theorem for samples from finite populations. *Studia Sci. Math. Hungar.*, **4**, 345–354 [in Russian].

Bloznelis, M. and Götze, F. (1997) An Edgeworth expansion for finite population *U*-statistics (1997). Preprint 97-012, SFB 343, Universität Bielefeld. Publication available on the Internet at http://www.mathematik.uni-bielefeld.de/sfb343/Welcome.html.

Bolthausen, E. and Götze, F. (1993) The rate of convergence for multivariate sampling statistics. *Ann. Statist.*, **21**, 1692–1710.

Callaert, H., Janssen, P. and Veraverbeke, N. (1980) An Edgeworth expansion for *U*-statistics. *Ann. Statist.*, **8**, 299–312.

Does, R.J.M.M. (1983) An Edgeworth expansion for simple linear rank statistics under the null-hypothesis. *Ann. Statist.*, **11**, 607–624.

Erdős, P. and Rényi, A. (1959) On the central limit theorem for samples from a finite population. *Publ. Math. Inst. Hungar. Acad. Sci.*, **4**, 49–61.

Götze, F. (1979) Asymptotic expansions for bivariate von Mises functionals. *Z. Wahrscheinlichkeitstheorie verw. Geb.*, **50**, 333–355.

Götze, F and van Zwet, W.R. (1991) Edgeworth expansions for asymptotically linear statistics. Preprint 91-034, SFB 343, Universität Bielefeld.

Helmers, R. and van Zwet, W.R. (1982) The Berry–Esseen bound for *U*-statistics. In S.S. Gupta and J.O. Berger (eds), *Statistical Decision Theory and Related Topics* III, Vol. 1, pp. 497–512. New York: Academic Press.

Hoeffding, W. (1963) Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.*, **58**, 13–30.

Höglund, T. (1978) Sampling from a finite population: a remainder term estimate. *Scand. J. Statist.*, **5**, 69–71.

Kokic, P.N. and Weber, N.C. (1990) An Edgeworth expansion for *U*-statistics based on samples from finite populations. *Ann. Probab.*, **18**, 390–404.

Kokic, P.N. and Weber, N.C. (1991) Rates of strong convergence for *U*-statistics in finite populations. *J. Austral. Math. Soc. Ser. A*, **50**, 468–480.

Nandi, H.K. and Sen, P.K. (1963) On the properties of *U*-statistics when the observations are not independent II. Unbiased estimation of the parameters of a finite population. *Calcutta Statist. Assoc. Bull.*, **12**, 124–148.

Prawitz, H. (1972) Limits for a distribution, if the characteristic function is given in a finite domain. *Skand. Aktuar. Tidskr.*, **5**, 138–154.

Robinson, J. (1978) An asymptotic expansion for samples from a finite population. *Ann. Statist.*, **6**, 1005–1011.

Serfling, R.J. (1980) *Approximation Theorems of Mathematical Statistics*. New York: Wiley.

Schneller, W. (1989) Edgeworth expansions for linear rank statistics. *Ann. Statist.*, **17**, 1103–1123.

van Zwet, W.R. (1982) On the Edgeworth expansion for the simple linear rank statistic. In B.V. Gnedenko, M.L. Puri and I. Vincze (eds), *Nonparametric Statistical Inference, Budapest 1980*, Vol. II, Colloq. Math. Soc. János Bolyai 32, pp. 889–909. Amsterdam: North-Holland.

van Zwet, W.R. (1984) A Berry–Esseen bound for symmetric statistics. *Z. Wahrscheinlichkeitstheorie verw. Geb.*, **66**, 425–440.

Zhao, L.C. and Chen, X.R. (1987) Berry–Esseen bounds for finite-population $U$-statistics. *Sci. Sinica Ser. A*, **30**, 113–127.

Zhao, L.C. and Chen, X.R. (1990) Normal approximation for finite-population $U$-statistics. *Acta Math. Appl. Sinica*, **6**, 263–272.