

THE BERRY-ESSEEN THEOREM FOR U -STATISTICS

BY HERMAN CALLAERT AND PAUL JANSSEN

Limburgs Universitair Centrum, Belgium

Assuming only the existence of the third absolute moment we prove that $\sup_x |P(\sigma_n^{-1}U_n \leq x) - \Phi(x)| \leq C\nu_3\sigma_g^{-3}n^{-1/2}$ where U_n is a U -statistic. This concludes a series of investigations on the Berry-Esseen theorem for U -statistics by Grams and Serfling, Bickel, and Chan and Wierman.

1. Introduction. Let $X_1, X_2, \dots, X_n, n \geq 2$, be i.i.d. random variables with common distribution function F . Define a U -statistic by $U_n = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} h(X_i, X_j)$ where h is a symmetric function of two variables with $Eh(X_1, X_2) = 0$ and such that $g(X_1) = E(h(X_1, X_2) | X_1)$ has a positive variance σ_g^2 . It then follows from Hoeffding (1948) that the distribution function (df) of $\sigma_n^{-1}U_n$ converges for $n \rightarrow \infty$ to the standard normal df Φ under the sole condition of the existence of $Eh^2(X_1, X_2)$. A study of the rate of this convergence started in 1973 with a paper by Grams and Serfling showing that $\sup_x |P(\sigma_n^{-1}U_n \leq x) - \Phi(x)|$ is of the order $O(n^{-r/(2r+1)})$, $n \rightarrow \infty$, when $Eh^{2r} < \infty$, leading to $O(n^{-1/2+\epsilon})$, $\epsilon > 0$, when h has finite moments of all orders. An order bound of exactly $O(n^{-1/2})$ was found by Bickel (1974) assuming U -statistics with bounded kernels h . Chan and Wierman (1977) succeeded in weakening considerably the assumptions of the previous theorems obtaining the order bound of $O(n^{-1/2})$ when the fourth moment exists and $O(n^{-1/2} \log^2 n)$ for kernels h with finite third absolute moment. We now prove that $O(n^{-1/2})$ can be attained requiring only the existence of the third absolute moment which is a natural assumption for a Berry-Esseen theorem. With only some easy additional computations it will be shown that $\sup_x |P(\sigma_n^{-1}U_n \leq x) - \Phi(x)| \leq C\nu_3\sigma_g^{-3}n^{-1/2}$ for all $n \geq 2$, indicating a nice analogy with the classical Berry-Esseen theorem.

The notation of this paper mainly adheres to that of Chan and Wierman and also the technique used in the first part of the proof is based on their work.

2. The Berry-Esseen bound. Let $\hat{U}_n = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} (g(X_i) + g(X_j))$ be the projection of U_n . Note that $E\hat{U}_n = 0$ and put $\hat{\sigma}_n^2 = E\hat{U}_n^2$ which equals $4\sigma_g^2n^{-1}$. Denote $(U_n - \hat{U}_n)\hat{\sigma}_n^{-1}$ by Δ_n and split it up into two parts Δ_n' and Δ_n'' such that $\Delta_n' = \binom{n}{2}^{-1}\hat{\sigma}_n^{-1} \sum_{1 \leq i < j \leq n} Y_{ij}$ with $Y_{ij} = h(X_i, X_j) - g(X_i) - g(X_j)$. The quantity c_n will be determined in the course of the proof and plays an essential role in obtaining various order bounds for several terms to be estimated.

THEOREM. Let $U_n = \binom{n}{2}^{-1} \sum_{1 \leq i < j \leq n} h(X_i, X_j)$ be a U -statistic such that $g(X_1) = E(h(X_1, X_2) | X_1)$ has positive variance σ_g^2 . If $\nu_3 = E|h(X_1, X_2)|^3 < \infty$ then there exists

Received March 1977; revised June 1977.

AMS 1970 subject classification. Primary 60F05.

Key words and phrases. Berry-Esseen bound, U -statistics.

an absolute constant C such that for all $n \geq 2$

$$(1) \quad \sup_x |P(\sigma_n^{-1}U_n \leq x) - \Phi(x)| \leq C\nu_3\sigma_g^{-3}n^{-\frac{1}{2}}.$$

PROOF. First note that

$$\sigma_g^2 = Eg^2(X_1) = Eh(X_1, X_2)h(X_1, X_3) \leq Eh^2(X_1, X_2) \leq \nu_3^{\frac{2}{3}} < \infty,$$

so that $\nu_3\sigma_g^{-3} \geq 1$. Also $E|g(X_1)|^3 \leq \nu_3$. Let $S_n = \hat{\sigma}_n^{-1}\hat{U}_n = n^{-\frac{1}{2}}\sigma_g^{-1} \sum_{i=1}^n g(X_i)$ which is a standardized sum of i.i.d. random variables with finite third absolute moment. Note that $\hat{\sigma}_n^{-1}U_n = (S_n + \Delta_n') + \Delta_n''$ and write

$$(2) \quad \begin{aligned} \sup_x |P(\hat{\sigma}_n^{-1}U_n \leq x) - \Phi(x)| \\ \leq \sup_x |P(S_n + \Delta_n' \leq x) - \Phi(x)| + P(|\Delta_n''| \geq a_n) + O(a_n). \end{aligned}$$

Then, with φ_x the characteristic function of X ,

$$(3) \quad \begin{aligned} \int_0^{\varepsilon n^{\frac{1}{2}}} t^{-1} |e^{-t^2/2} - \varphi_{S_n + \Delta_n'}(t)| dt \\ \leq \int_0^{\varepsilon n^{\frac{1}{2}}} t^{-1} |e^{-t^2/2} - \varphi_{S_n}(t)| dt + \int_0^{\varepsilon n^{\frac{1}{2}}} t^{-1} |\varphi_{S_n}(t) - \varphi_{S_n + \Delta_n'}(t)| dt. \end{aligned}$$

From the proof of the classical Berry–Esseen theorem it follows that

$$(4) \quad \int_0^{\varepsilon n^{\frac{1}{2}}} t^{-1} |e^{-t^2/2} - \varphi_{S_n}(t)| dt \leq C_1 E|g(X_1)|^3 \sigma_g^{-3} n^{-\frac{1}{2}} \leq C_1 \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}}$$

for $\varepsilon = \sigma_g^3/E|g(X_1)|^3$. Throughout the C_k are absolute constants.

We now start the estimation of the last integral in (3). Writing η for the characteristic function $\varphi_{g(X_1)}$ we have $|\eta(\vartheta)| \leq \exp(-\frac{1}{3}\vartheta^2\sigma_g^2)$ for $|\vartheta| \leq \varepsilon\sigma_g^{-1}$ and ε as above. Also note that $E(f(X_i)Y_{ij}) = 0$ for any bounded Borel-measurable function f . Then $|\varphi_{S_n}(t) - \varphi_{S_n + \Delta_n'}(t)| = |Ee^{itS_n}(1 - e^{it\Delta_n'})| \leq |Ee^{itS_n}it\Delta_n'| + \frac{1}{2}t^2E\Delta_n'^2$. Now, for $0 \leq t \leq \varepsilon n^{\frac{1}{2}}$ and $n \geq 4$,

$$\begin{aligned} |Ee^{itS_n}\Delta_n'| &= \frac{1}{2}\sigma_g^{-1}n^{\frac{1}{2}}\binom{n}{2}^{-1} \sum_{1 \leq i < j \leq c_n} E(\exp[itn^{-\frac{1}{2}}\sigma_g^{-1} \sum_{k \neq i, j} g(X_k)]) \\ &\quad \times E(\exp[itn^{-\frac{1}{2}}\sigma_g^{-1}[g(X_i) + g(X_j)]]Y_{ij}) \\ &= \frac{1}{2}\sigma_g^{-1}n^{\frac{1}{2}}\binom{n}{2}^{-1} |\eta(n^{-\frac{1}{2}}\sigma_g^{-1}t)|^{n-2} \\ &\quad \times |\sum_{1 \leq i < j \leq c_n} E(\exp[itn^{-\frac{1}{2}}\sigma_g^{-1}g(X_i)] - 1)(\exp[itn^{-\frac{1}{2}}\sigma_g^{-1}g(X_j)] - 1)Y_{ij}| \\ &\leq \frac{1}{2}\sigma_g^{-3}n^{-\frac{1}{2}}t^2e^{-t^2/6} E|g(X_i)g(X_j)Y_{ij}| \\ &\leq \frac{3}{2}\nu_3\sigma_g^{-3}n^{-\frac{1}{2}}t^2e^{-t^2/6}. \end{aligned}$$

Hence, for $d_n \leq \varepsilon n^{\frac{1}{2}}$,

$$(5) \quad \begin{aligned} \int_0^{d_n} t^{-1} |Ee^{itS_n}it\Delta_n'| dt &\leq \frac{3}{2}\nu_3\sigma_g^{-3}n^{-\frac{1}{2}} \int_0^{d_n} t^2e^{-t^2/6} dt \\ &\leq \frac{9}{4}(6\pi)^{\frac{1}{2}}\nu_3\sigma_g^{-3}n^{-\frac{1}{2}}. \end{aligned}$$

Remark that the last estimate, which is also valid for $n = 2, 3$ as may be seen by direct computations, is independent of the choice of d_n . Also c_n does not play any particular role at this moment. Further, since $E\Delta_n'^2 \leq Eh^2(X_1, X_2)\sigma_g^{-2}n^{-1} \leq \nu_3\sigma_g^{-3}n^{-1}$, we have

$$(6) \quad \frac{1}{2}E\Delta_n'^2 \int_0^{d_n} t dt \leq \frac{1}{4}\nu_3\sigma_g^{-3}n^{-1}d_n^2$$

which for $d_n = n^{\frac{1}{2}}$ yields the upper bound $\frac{1}{4}\nu_3\sigma_g^{-3}n^{-\frac{1}{2}}$. The estimates (5) and (6)

are sufficient for all n such that $\varepsilon n^{\frac{1}{2}} \leq n^{\frac{1}{2}}$. If $d_n < \varepsilon n^{\frac{1}{2}}$ we write

$$\begin{aligned} & |E e^{itS_n}(1 - e^{it\Delta'_n})| \\ & \leq |E \exp[itn^{-\frac{1}{2}}\sigma_g^{-1} \sum_{k>c_n} g(X_k)]| |E \exp[itn^{-\frac{1}{2}}\sigma_g^{-1} \sum_{k\leq c_n} g(X_k)]| (1 - e^{it\Delta'_n})| \\ & \leq tE|\Delta'_n| |\eta(n^{-\frac{1}{2}}\sigma_g^{-1}t)|^{n-c_n} \\ & \leq \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}} \exp(-(n - c_n)n^{-1}t^2/3). \end{aligned}$$

Hence

$$(7) \quad \int_{d_n}^{\varepsilon n^{\frac{1}{2}}} t^{-1} |\varphi_{S_n}(t) - \varphi_{S_n + \Delta'_n}(t)| dt \leq \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}} \int_{d_n}^{\varepsilon n^{\frac{1}{2}}} \exp(-(n - c_n)n^{-1}t^2/3) dt.$$

Here the flexibility of a choice for d_n and c_n already becomes clear. As an example we could take $d_n = n^{\frac{1}{2}}$ and $c_n = [n - 3n^{\frac{1}{2}} \log n]$ which is, in the Chan and Wierman paper, a crucial choice for obtaining the overall order bound of $O(n^{-\frac{1}{2}} \log^{\frac{1}{2}} n)$. The estimate in (7) then becomes $C_2 \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}}$ but it is easily seen that many other d_n and c_n provide an analogous bound. From the classical Berry-Esseen argument, together with (4), (5), (6), (7) and a suitable choice of d_n and c_n we obtain

$$\sup_x |P(S_n + \Delta'_n \leq x) - \Phi(x)| \leq C_3 \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}}$$

yielding the desired result for the first term on the right-hand side of (2). For an estimate of the other terms in (2) we use the following

LEMMA. *With the notations and the assumptions as in the above theorem one has*

$$(8) \quad E|\Delta_n''|^3 \leq C_4 \nu_3 \sigma_g^{-3} (n - c_n)^{\frac{3}{2}} n^{-3} \quad \text{for } n = 2, 3, \dots$$

We postpone the proof of the lemma to the end of this paper.

It now follows from the Markov inequality that

$$(9) \quad P(|\Delta_n''| \geq a_n) \leq C_4 \nu_3 \sigma_g^{-3} n^{-3} (n - c_n)^{\frac{3}{2}} a_n^{-3}$$

and again we have a lot of freedom in choosing a_n and c_n . One systematic way consists in taking

$$(10) \quad a_n = [n^{-3}(n - c_n)^{\frac{3}{2}}]^{\frac{1}{2}}$$

and then choosing c_n such that $a_n \leq C_5 n^{-\frac{1}{2}}$. This yields

$$(11) \quad P(|\Delta_n''| \geq \nu_3 \sigma_g^{-3} a_n) \leq P(|\Delta_n''| \geq a_n) \leq C_4 \nu_3 \sigma_g^{-3} a_n$$

which is sufficient for obtaining the estimate $C_6 \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}}$ for the last two terms in (2). We finally note that $\sigma_n^2 = \binom{n}{2}^{-1} (Eh^2(X_1, X_2) + 2(n - 2)\sigma_g^2)$, and hence for $n \geq 2$

$$\left| \frac{\sigma_n}{\hat{\sigma}_n} - 1 \right| \leq \left| \frac{\sigma_n^2}{\hat{\sigma}_n^2} - 1 \right| \leq \frac{1}{n - 1} \left(1 + \frac{Eh^2(X_1, X_2)}{2\sigma_g^2} \right) \leq 3\nu_3 \sigma_g^{-3} n^{-1}$$

completing the proof of the theorem.

3. Some remarks on the quantities a_n , c_n , and d_n . If one only wants the order bound $O(n^{-\frac{1}{2}})$ in the Berry-Esseen theorem, the above proof could be somewhat

simplified by making in advance a suitable choice of a_n , c_n and d_n and then, without further discussion, showing that (1) holds. We preferred not to follow this way for the following reason. If one is interested in some numerical value for the constant C in (1) the possibility of altering a_n , c_n and d_n within some suitable ranges may provide easier and sharper numerical bounds on the various C_k appearing in the proof. In fact, within our framework, one is able to distinguish between dominant terms yielding $O(n^{-\frac{1}{2}})$ and terms which are $o(n^{-\frac{1}{2}})$ for $n \rightarrow \infty$.

To make this explicit we determine the ranges for the quantities under consideration. From (7) and (10) it follows that $n - c_n$ may vary from $O(n^{\frac{1}{2}} \log n)$ to $O(n^{\frac{3}{2}})$ while d_n is at most $O(n^{\frac{1}{2}})$ by (6). Taking $n - c_n = n^{\frac{3}{2}}$ in (7) we find $O(n^{\frac{1}{2}}(\log n)^{\frac{1}{2}})$ as a lower bound for d_n . Finally a_n is related to c_n by (10). It now becomes clear that apart from the estimate in (4) which obviously cannot be smaller than $O(n^{-\frac{1}{2}})$, the only dominant term appears in (5). All other estimates relevant in the proof of (1) can be made $o(n^{-\frac{1}{2}})$ by staying away from the end-points of the ranges indicated above. This makes it possible to write (1) as

$$\sup_x |P(\sigma_n^{-1}U_n \leq x) - \Phi(x)| \leq C_1' \nu_3 \sigma_g^{-3} n^{-\frac{1}{2}} + C_2' o(n^{-\frac{1}{2}})$$

where C_1' is considerably smaller than C_2' .

PROOF OF THE LEMMA. Define ξ_j by

$$(12) \quad \binom{n}{2} \hat{\sigma}_n \Delta_n'' = \sum_{j=c_n+1}^n \sum_{i=1}^{j-1} Y_{ij} = \sum_{j=c_n+1}^n \xi_j.$$

We have $\xi_1 = 0$ and $E(\xi_{j+1} | \xi_1, \dots, \xi_j) = 0$ a.s. for $j = 1, 2, \dots$. Hence the ξ_j are martingale summands and, by optional skipping, $V_k = \sum_{j=c_n+1}^{c_n+k} \xi_j$ forms a martingale, $k = 1, 2, \dots, n - c_n$. Applying a theorem of Dharmadhikari, Fabian and Jogdeo (1968) we get for $k = n - c_n$

$$(13) \quad E|V_{n-c_n}|^3 \leq 2^{12}(n - c_n)^{\frac{3}{2}} \max_{c_n+1 \leq j \leq n} E|\xi_j|^3.$$

But for fixed $j \geq c_n + 1$, $W_k = \sum_{i=1}^k Y_{ij}$, $k = 1, 2, \dots, j - 1$, is also a martingale and the same argument yields for $j = 2, 3, \dots$

$$(14) \quad E|\xi_j|^3 = E|W_{j-1}|^3 \leq 2^{12}(j - 1)^{\frac{3}{2}} \max_{1 \leq i \leq j-1} E|Y_{ij}|^3 \leq 2^{12}3^3(j - 1)^{\frac{3}{2}}\nu_3.$$

Then, from (12), (13) and (14)

$$E|\binom{n}{2} \hat{\sigma}_n \Delta_n''|^3 \leq 2^{24}3^3(n - c_n)^{\frac{3}{2}}(n - 1)^{\frac{3}{2}}\nu_3$$

and hence

$$E|\Delta_n''|^3 \leq 2^{24}3^4(n - c_n)^{\frac{3}{2}}n^{-3}\nu_3\sigma_g^{-3} \quad \text{for } n = 2, 3, \dots,$$

which proves the lemma.

4. The c -sample case. Although we proved the Berry-Esseen theorem only for a one-sample U -statistic of order two, the result remains valid for the general case of multisample U -statistics of arbitrary order provided that the minimum sample size tends to infinity. The proof is based on the same ideas as used in this paper but becomes computationally more involved and will not be given

here. For more details we refer to the Ph. D. thesis to be completed by Paul Janssen.

NOTE. In the proof of the theorem the quantity Δ_n'' has been separated from the U -statistic and handled by the Markov inequality. An alternative procedure consists in writing $\Delta_n = \Delta_{n_1} + \Delta_{n_2}$, then making a Taylor-expansion of $Ee^{it\Delta_{n_2}}$ and using an independence argument. This method will be displayed in a forthcoming paper on the Edgeworth expansion for U -statistics, co-authored by N. Veraverbeke and the authors.

Acknowledgment. The authors thank Professor W. R. van Zwet for the suggestion of showing the dependence of the order bound on F and h explicitly.

REFERENCES

- [1] BICKEL, P. J. (1974). Edgeworth expansions in nonparametric statistics. *Ann. Statist.* **2** 1-20.
- [2] CALLAERT, H., JANSSEN, P. and VERAVERBEKE, N. (1977). Edgeworth expansion for U -statistics. Unpublished manuscript.
- [3] CHAN, Y. and WIERMAN, J. (1977). On the Berry-Esseen theorem for U -statistics. *Ann. Probability* **5** 136-139.
- [4] DHARMADHIKARI, S. W., FABIAN, V. and JOGDEO, K. (1968). Bounds on the moments of martingales. *Ann. Math. Statist.* **39** 1719-1723.
- [5] GRAMS, W. F. and SERFLING, R. J. (1973). Convergence rates for U -statistics and related statistics. *Ann. Statist.* **1** 153-160.
- [6] Hoeffding, W. (1948). A class of statistics with asymptotically normal distribution. *Ann. Math. Statist.* **19** 293-325.

DEPARTMENT OF MATHEMATICS
LIMBURGS UNIVERSITAIR CENTRUM
B-3610 DIEPENBEEK, BELGIUM