Asymptotic Distributions of Weighted U-Statistics of Degree 2
Author(s): Kevin A. O'Neil and Richard A. Redner
Source: *The Annals of Probability*, Vol. 21, No. 2 (Apr., 1993), pp. 1159-1169
Published by: Institute of Mathematical Statistics
Stable URL: http://www.jstor.org/stable/2244694
Accessed: 04-12-2015 13:37 UTC

http://www.jstor.org

# ASYMPTOTIC DISTRIBUTIONS OF WEIGHTED U-STATISTICS OF DEGREE 2

By Kevin A. O'Neil and Richard A. Redner[1]

*University of Tulsa*

The limiting distribution of weighted U-statistics of degree 2 is found for a wide class of weights, including uniform weights. Nonnormal limits can occur for both degenerate and nondegenerate kernels. A compact expression is given for the cumulants of the distribution. Incomplete and randomly weighted U-statistics are also analyzed.

**1. Introduction.** Let $X_1, X_2, \ldots$ be iid random variables with distribution $F$, and $f(x, y)$ be a real symmetric kernel function with mean zero and finite variance: $E_F(f(X_1, X_2)) = 0$, $E_F(f^2(X_1, X_2)) < \infty$. Given a bounded symmetric "weight function" $a: \mathbf{N}^2 \to \mathbf{R}$, consider the sequence

$$(1.1) \qquad U_n = \sum_{1 \le i < j \le n} a(i, j) f(X_i, X_j).$$

When $a \equiv 1$, these sums are known as U-statistics of degree 2; the recent book by Lee (1990) collects much of the literature on the subject. Of particular interest is the asymptotic distribution of these statistics. If the kernel $f$ is nondegenerate, that is, if the function $\mu(x) = E_F(f(x, X_1))$ has positive variance, the limiting distribution is normal [Hoeffding (1948)]. For degenerate kernels, the limiting distribution of $U_n$ was described by Gregory (1977) and Serfling (1980) in terms of the eigenvalues of an operator constructed from $f$, as a sum of independent chi-square variables; see also Rubin and Vitale (1980) and Dynkin and Mandelbaum (1983).

When the weight function $a$ is nonconstant, $U_n$ is called a weighted U-statistic. Conditions on $a$ sufficient to produce asymptotic normality when $f$ is nondegenerate were reported by Shapiro and Hubert (1979). An important special case is the reduced or incomplete U-statistic, where $a$ takes only the values 0 and 1. Designs (weight functions) resulting in both normal and nonnormal convergence have been analyzed by Blom (1976), Brown and Kildea (1978) and Janson (1984). However, general weighted U-statistics with nonnormal asymptotic distributions have not, to our knowledge, been described in the literature. Such U-statistics are the subject of this paper.

We find that nonnormal limits can occur when the kernel is degenerate, and when the kernel is nondegenerate but the weights are degenerate in a certain sense. We give a compact expression for the cumulants of the asymptotic distribution, involving certain expectations of both the kernel and the weights.

1159

Knowledge of the cumulants makes nonnormality of the limit transparent. A theorem establishing normal convergence for a wide class of incomplete $U$-statistics with degenerate kernels is given, extending a theorem of Brown and Kildea (1978).

We also investigate several weighting schemes for which the preceding expectations can be computed easily. One simple scheme, prominent in applications in physics, is a factored form: $a(i, j) = e_i e_j$. The interesting conclusion is that, up to a scaling constant, the limiting distribution is indifferent to the choice of "charges" $e_i$. A generalization where the variables $X_i$ are assigned to groups and the weights are determined by the groups is discussed. Randomly assigned weights are also analyzed.

The limiting distributions are found by the method of moments: The limiting moments are found and shown to determine a distribution. Although the moment computations can be generalized to weighted $U$-statistics of degree greater than 2 [cf. O'Neil and Redner (1992)], these moments are in general insufficient to determine a distribution. The existence of limiting distributions for higher order weighted $U$-statistics is an open question.

The organization of the paper is as follows. In Sections 2 and 3 the main theorems (Theorems 2.1 and 3.1) are proved, handling degenerate and nondegenerate kernels, respectively. The case of randomly assigned weights is treated in Corollary 2.2. Section 4 deals with incomplete $U$-statistics (Theorem 4.1), and the last section discusses some simple weight functions which satisfy the hypotheses of the theorems. Examples from and connections with the literature are given throughout.

**2. Weighted $U$-statistics with degenerate kernels.** Let $f$, $a$, $X_i$ and $U_n$ be defined as before. (One may also consider a sequence of kernels $f_n$ and weight functions $a_n$; a trivial modification of the proofs that follow handles this case.) Our first theorem discusses the asymptotics of $U_n$ when the kernel $f$ is degenerate; the limiting cumulants are expressed in terms of expectations over "cycles":

$$(2.1) \qquad I_k = E_F\big( f(X_1, X_2) f(X_2, X_3) \cdots f(X_k, X_1)\big).$$

(These expectations are finite since $f$ has finite variance; see Remark 1.) The $I_k$ have been shown to determine the limit of unweighted $U$-statistics [O'Neil and Redner (1991)].

THEOREM 2.1. *Let f be degenerate, and define $I_k$ as before. Suppose that for all $k \geq 2$, the following limits exist:*

$$(2.2) \qquad w_k = \lim_{n \to \infty} n^{-k} \sum a(i_1, i_2) a(i_2, i_3) \cdots a(i_k, i_1),$$

*where the sum is over all $k$-tuples $(i_1, \ldots, i_k)$ of distinct integers from 1 to $n$. If $w_2 > 0$, then the sequence $n^{-1} U_n$ converges in distribution to a limit with $k$th*

*cumulant* $(k - 1)! I_k w_k / 2$ *for* $k \geq 2$, *that is, with moment generating function*

(2.3)
$$\phi(t) = \exp\left[\sum_{2}^{\infty} \frac{I_k w_k}{2k} t^k\right].$$

PROOF. The proof is by the method of moments: We show that the moments of $n^{-1}U_n$ tend to limits corresponding to the previously claimed cumulants, and that these moments determine the distribution. First we assume that $f$ is bounded; this restriction will be lifted later.

Begin by evaluating the $m$th moment of $n^{-1}U_n$:

(2.4)
$$E_F\left(\left(n^{-1} \sum a(i, j) f(X_i, X_j)\right)^m\right)$$
$$= n^{-m} \sum a(i_1, j_1) \cdots a(i_m, j_m) E_F\left(f(X_{i_1}, X_{j_1}) \cdots f(X_{i_m}, X_{j_m})\right),$$

where the summation is over all ordered $m$-tuples $((i_1, j_1), \ldots, (i_m, j_m))$ of pairs of integers from 1 to $n$ with $i_k < j_k$ for all $k$.

Since $f$ is degenerate, the preceding expectations, in which any integer appears only once in a subscript, all vanish. Furthermore, those terms in which the $m$-tuple of pairs contains fewer than $m$ distinct integers in the subscripts make only an $O(1/n)$ contribution to the moment: $a$ and $f$ are bounded and there are only $O(n^{(m-1)})$ such terms.

Now consider the terms in which the $m$-tuple of pairs contains exactly $m$ distinct integers in subscripts, each appearing exactly twice. If the $m$ pairs are thought of as vertices, and the $m$ distinct integers are used to label edges connecting these vertices, then all these terms can be represented as a collection of cycles. The expectations in (2.4) then take the form $I_{n_1}^{e_1} \cdots I_{n_r}^{e_r}$ for some positive integers $n_1, \ldots, n_r, e_1, \ldots, e_r$, $1 < n_1 < \cdots < n_r \leq m$, satisfying the relation $n_1 e_1 + \cdots + n_r e_r = m$. Ignoring $O(1/n)$ terms, the $m$th moment of $n^{-1}U_n$ is

(2.5)
$$\sum_{[m]} I_{n_1}^{e_1} \cdots I_{n_r}^{e_r} \cdot n^{-m} \sum a(i_1, j_1) \cdots a(i_m, j_m),$$

where $[m]$ indicates that the first summation is over $n_1, \ldots, n_r, e_1, \ldots, e_r$ as before, and the second summation is over all the $m$-tuples of pairs $((i_1, j_1), \ldots, (i_m, j_m))$ in which exactly $m$ distinct integers appear, $1 \leq i_k < j_k \leq n$ for all $k$, and the graph of the $m$-tuple consists of $e_k$ cycles of order $n_k$, $k = 1, \ldots, r$.

The second summation can be simplified in the following way. Fix the graph $G$ (that is, $n_1, \ldots, n_r, e_1, \ldots, e_r$), and note that the cycles and then the vertices of $G$ can be ordered so that the first $n_1$ vertices belong to the first cycle of order $n_1$, the second $n_1$ vertices belong to the second one (if $e_1 > 1$) and so forth; the last $n_r$ vertices make up the last cycle of order $n_r$. Now consider the sum

(2.6)
$$S(G, n) = n^{-m} \sum a(i_1, j_1) \cdots a(i_m, j_m),$$

where the summation is over all sets of $m$ distinct integers from 1 to $n$ which

appear in the summand in the way determined by this ordering of $G$. That is, $j_1 = i_2, j_2 = i_3, \ldots, j_{n_1} = i_1$ and so forth.

The $m$ factors of the terms of $S(G, n)$ may be permuted to form all the terms of the second summation of (2.5). However, it is not correct to replace this summation by $m! S(G, n)$, for the terms obtained in this way are not distinct. Each cycle of order $k > 2$ has a symmetry group of order $2k$ that results in duplications; similar overrepresentation is made when $k = 2$. Moreover, cycles of the same order are indistinguishable in $G$ but distinguished when $G$ is ordered to compute $S(G, n)$; this results in overrepresentation by a factor of $e_1! \ldots e_r!$. Taking these factors into account, we find that (2.5) becomes

$$(2.7) \qquad m! \sum_{[m]} \frac{\left(I_{n_1}/2n_1\right)^{e_1}}{e_1!} \cdots \frac{\left(I_{n_r}/2n_r\right)^{e_r}}{e_r!} \cdot S(G, n).$$

We now observe that $S(G, n) \to w_{n_1}^{e_1} \cdots w_{n_r}^{e_r}$ as $n \to \infty$. Simply replace $w_{n_1}$ by $n^{-n_1} \Sigma a(i_1, i_2) \cdots a(i_{n_1}, i_1)$ and so on, and expand the product. Those terms which do not appear in $S(G, n)$ are all products of the form $n^{-m} a(i_1, j_1) \cdots a(i_m, j_m)$, where the integers $i_1, j_1, \ldots, i_m, j_m$ take on fewer than $m$ distinct values. There are $O(n^{m-1})$ of these terms, and thus their sum has limit 0.

The limiting value of the $m$th moment of $n^{-1} U_n$ has now been established:

$$(2.8) \qquad \mu_m = m! \sum_{[m]} \frac{\left(I_{n_1} w_{n_1}/2n_1\right)^{e_1}}{e_1!} \cdots \frac{\left(I_{n_r} w_{n_r}/2n_r\right)^{e_r}}{e_r!}.$$

Next we note the formal relation

$$(2.9) \qquad \sum_0^\infty \frac{\mu_m}{m!} t^m = \exp\left[ \sum_2^\infty \frac{I_k w_k}{2k} t^k \right].$$

To see this, rearrange the right-hand side,

$$(2.10) \quad \exp\left[ \sum_2^\infty \frac{I_k w_k}{2k} t^k \right] = \prod_2^\infty \exp\left( \frac{I_k w_k t^k}{2k} \right) = \prod_{k=2}^\infty \sum_{j=0}^\infty \frac{(I_k w_k/2k)^j}{j!} t^{jk},$$

expand and then compare terms.

We show now that the cumulant generating function [in brackets on the right-hand side of (2.9)] has positive radius of convergence; (2.10) then shows that the moment generating function has as well, and as a consequence the moments determine a distribution [Feller (1968)]. Clearly $w_k \le M^k$, where $M$ is a bound for $a$, and repeated application of the Schwarz inequality gives the relation $I_k \le I_2^{k/2}$. Hence the series in (2.10) converge for $t^2 < 1/(M^2 I_2)$.

The final step in the proof is to remove the restriction that $f$ is bounded. Let $X$ and $Y$ be independent random variables with distribution $F$. We now assume only that $E_F(f^2(X, Y))$ is finite. Given a positive integer $M$, define a "mollified kernel" $\hat{f}_M(x, y)$ by setting $\hat{f}_M(x, y) = f(x, y)$ if $|f(x, y)| \le M$ and

zero otherwise. Then define the projection onto its degenerate part:

$$g_M(x, y) = f_M(x, y) - E_F(f_M(x, Y)) - E_F(f_M(X, y)) + E_F(f_M(X, Y)).$$

Clearly, $g_M$ is bounded and $g_M \to f$ a.e. as $M \to \infty$. Now we write

$$n^{-1} \sum a(i, j) f(X_i, X_j) = S_n(g_M) + S_n(f - g_M)$$

(2.11)
$$= n^{-1} \sum a(i, j) g_M(X_i, X_j)$$
$$+ n^{-1} \sum a(i, j)(f - g_M)(X_i, X_j).$$

It is easy to show that the variance of $S_n(f - g_M)$ goes to zero uniformly in $n$ as $M \to \infty$, and for each $M$, we have proved that $S_n(g_M)$ tends to a limiting distribution as $n \to \infty$. It suffices then to show that for each $k > 1$, the expectation

(2.12)          $$I_k(M) = E_F(g_M(X_1, X_2) \cdots g_M(X_k, X_1))$$

tends to $I_k$ as $M \to \infty$. The integrand of $I_k(M)$ converges to the integrand of $I_k$ a.e., so the desired limit follows from the dominated convergence theorem if the integrand can be shown to be dominated by an integrable function. To this end, define ·

(2.13)   $$F(x, y) = |f(x, y)| + E_F(|f(x, Y)|) + E_F(|f(X, y)|) + 1.$$

Clearly $|g_M| \le F$ for large $M$, so that the integrand of $I_k(M)$ is dominated by

(2.14)                    $$F(x_1, x_2) \cdots F(x_k, x_1),$$

which is integrable (again by the representation of Remark 1) if $F^2$ is integrable. Verification of this last fact amounts to checking that $E_F(|f(x, Y)|)$ has finite variance. This completes the proof. $\square$

REMARK 1.   The (unweighted) $U$-statistic is a special case, with $w_k = 1$ for all $k$, so Theorem 2.1 is valid for ordinary $U$-statistics as well. Indeed, the formula follows from the expression for the characteristic function given in Gregory (1977) and Serfling (1980) when the $I_k$ are related to the representation of $f$ used there. Specifically we have $f(x, y) = \sum \lambda_i \phi_i(x) \phi_i(y)$ for certain functions $\phi_i$ which satisfy the relations $E_F(\phi_i(X) \phi_j(X)) = \delta_{ij}$. From this representation it follows that $I_k = \sum \lambda_i^k$. The sums converge because, by hypothesis, $I_2$ converges.

REMARK 2.   From the formula for the cumulant generating function one sees that the limit distribution is normal only if $I_k w_k = 0$ for all $k > 2$. This never holds for the unweighted case, because the sums $\sum \lambda_i^{2k} = I_{2k}$ are positive. However, normality can occur when a sequence of degenerate kernels is used, by letting the support of the kernel shrink; see Jammalamadaka and Janson (1986).

REMARK 3.   The limiting distribution is symmetric iff all odd cumulants vanish. For unweighted $U$-statistics, this is a condition on the kernel: the limit

is symmetric iff the distribution of eigenvalues $\lambda_i$ is symmetric about the origin. For weighted $U$-statistics, symmetry may be a consequence of the weights. For example, the weight function $a(i, j) = 1 - (-1)^{i+j}$ yields limits $w_{2k+1} = 0$ for all $k > 0$.

REMARK 4.   There is a canonical expression for (unweighted) $U$-statistics of higher degree as the sum of a hierarchy of $U$-statistics having increasing degree and decreasing variance; see for example Lee (1990). The asymptotic behavior therefore depends only on the term with largest variance. Theorem 2.1 describes the limiting distribution if that term has degree 2. On the other hand, the decomposition of weighted $U$-statistics does not in general give a sequence with decreasing variance, so that these projections are less useful. The next section provides an explicit example of this.

Theorem 2.1 is easily extended to $U$-statistics with randomly assigned weights. If $S$ is a set of positive integers, let $G(S) = \{a(i, j)|i \text{ and } j \in S\}$. We consider weights $a(i, j)$ that are bounded random variables, independent of all the $X_i$, and with the additional independence condition: if $S_1$ and $S_2$ are disjoint, then $G(S_1)$ and $G(S_2)$ are independent.

COROLLARY 2.2.   Let $f$, $X_i$ and $U_n$ be as in Theorem 2.1, and $a(i, j)$ be as before. Suppose that the limits

$$(2.15) \qquad \overline{w}_k = \lim_{n \to \infty} n^{-k} \sum E\big(a(i_1, i_2) \cdots a(i_k, i_1)\big)$$

exist for $k \geq 2$ and that $\overline{w}_2 > 0$. Then $n^{-1}U_n$ converges in distribution to a limit with $k$th cumulant $(k - 1)! I_k \overline{w}_k / 2$.

PROOF.   In the expression for the $m$th moment of $n^{-1}U_n$ in (2.4), replace products of weights by the expectation of the product. The proof then proceeds as in Theorem 2.1. The independence condition is needed to factor the expectation of the product of weights into the product of the expectations over the connected components of the graph of that product. □

As a simple example, take $a(i, j)$ to be independent random variables with means $\mu_{ij}$ and variance $\sigma_{ij}^2$. Suppose that $\sigma^2 = \lim n^{-2} \sum \sigma_{ij}^2$ exists and $n^{-1} \sum \mu_{ij} f(X_i, X_j)$ converges in distribution to a limit $Y$ by Theorem 2.1; that is, $w_2, w_3, \ldots$ exist. A short calculation shows that $\overline{w}_2 = \sigma^2 + w_2$ and $\overline{w}_k = w_k$ for all $k > 2$. Thus the cumulant generating function splits and the limiting distribution of $n^{-1}U_n$ is that of the sum of two independent random variables, one with distribution $N(0, I_2\sigma^2/2)$ and the other with distribution $Y$. For a related result, see Janson (1984).

**3. Nondegenerate kernels.**   In this section we suppose that $f$ is a nondegenerate kernel with zero mean and finite variance. The type of limiting distribution of $U_n$ now depends on the weights. Recall that the projection of $f$

onto its degenerate part $\tilde{f}$ involves the conditional expectations $h(x) = E_F(f(x, Y))$; specifically, $\tilde{f}(x, y) = f(x, y) - h(x) - h(y)$. We can define a similar conditional expectation for the weights for each $n$:

$$(3.1) \qquad \nu_i(n) = \langle a(i, \cdot) \rangle_n = \frac{1}{n} \sum_{j=1}^{n} a(i, j)$$

[taking $a(i, i) = 0$]. Likewise we write $\langle \nu_i^2(n) \rangle_n = n^{-1} \Sigma \nu_i^2(n)$, and so forth.

In the proof of Theorem 2.1, it was found that the moments were determined by the expectations over cycles $I_k$ and $w_k$. In the nondegenerate case we will need the expectation of $\tilde{f}$ over cycles,

$$(3.2) \qquad \tilde{I}_k = E_F\big( \tilde{f}(X_1, X_2) \tilde{f}(X_2, X_3) \cdots \tilde{f}(X_k, X_1) \big),$$

as well as expectations over "chains,"

$$(3.3) \qquad \begin{aligned} J_k &= E_F\big( h(X_1) \tilde{f}(X_1, X_2) \cdots \tilde{f}(X_{k-2}, X_{k-1}) h(X_{k-1}) \big), \\ z_k &= \lim_{n \to \infty} n^{-k} \sum a(i_0, i_1) a(i_1, i_2) \cdots a(i_{k-2}, i_{k-1}) a(i_{k-1}, i_k), \end{aligned}$$

where the sums are taken over all $(k + 1)$-tuples of distinct integers from 1 to $n$.

THEOREM 3.1. *Let f be a nondegenerate kernel with zero mean.*

(i) *If $n \langle \nu_i^2(n) \rangle_n \to \infty$, then $U_n$ is asymptotically normal.*
(ii) *If $\langle \nu_i^2(n) \rangle_n$ is $o(1/n)$, the limits $w_2, w_3, \ldots$ all exist and $w_2 > 0$, then $n^{-1} U_n$ has the same limiting distribution as $n^{-1} \Sigma_{i < j} a(i, j) \tilde{f}(X_i, X_j)$.*
(iii) *Suppose that the sums $v_i(n) = n^{-1/2} \Sigma_1^n a(i, j)$ are bounded [taking $a(i, i) = 0$], so that $\langle \nu_i^2(n) \rangle_n$ is $O(1/n)$. If the limits $w_k$ and $z_k$ exist for $k > 1$ and $z_2 > 0$, then the sequence $n^{-1} U_n$ tends to a limiting distribution with kth cumulant $((k - 1)! \tilde{I}_k w_k + k! J_k z_k)/2$; that is, with moment generating function*

$$(3.4) \qquad \phi(t) = \exp\left[ \left(\frac{1}{2}\right) \sum_{k=2}^{\infty} \left( \frac{\tilde{I}_k w_k}{k} + J_k z_k \right) t^k \right].$$

PROOF. Using the definition of $\tilde{f}$, we have

$$(3.5) \quad \sum_{j<k} a(j, k) f(X_j, X_k) = \sum_{j<k} a(j, k) \tilde{f}(X_j, X_k) + n \sum h(X_i) \nu_i(n).$$

The variance of the first sum on the right is $O(n^2)$, while that of the second is $O(n^2) \cdot n \langle \nu_i^2(n) \rangle_n$. If $n \langle \nu_i^2(n) \rangle_n \to \infty$, then the second sum dominates and the central limit theorem applies, but if $\langle \nu_i^2(n) \rangle_n$ is $o(1/n)$, the first sum dominates and Theorem 2.1 applies. This establishes (i) and (ii).

The proof of (iii) follows the pattern of the proof of Theorem 2.1 closely and is only sketched here. We may assume that $f$ is bounded and write

$$(3.6) \qquad n^{-1}U_n = n^{-1} \sum_{j<k} a(j,k)\tilde{f}(X_j, X_k) + n^{-1/2} \sum h(X_i)v_i(n).$$

Begin by noting that in the expansion of $\langle(n^{-1}U_n)^m\rangle$, boundedness of $a$ and $v_i(n)$ implies that the limit as $n \to \infty$ is determined by those terms in which exactly $m$ distinct integers appear as subscripts, each appearing exactly twice. In particular, an even number of the $m$ factors must be $h$'s. Each term can be described graphically, with each factor of $\tilde{f}$ represented by a vertex touching two edges and each factor of $h$ by a vertex touched by one edge. Each connected component of each graph is either a cycle or a chain, and the expectations therefore can be expressed as products of powers of $\tilde{I}_k$ and $J_k$. When expressing the graph in standard form with a standard ordering of the vertices, one must note the symmetry group of a chain has order 2. The formula for the moments, and hence for the cumulants, now follows as in Theorem 2.1. Finally it is easily checked that the limiting distribution is determined by the limits of the moments. This completes the proof. $\square$

A simple example is provided by the weight function $a(i,j) = (-1)^{i+j}$, for which $\langle v_i^2(n)\rangle_n$ is $O(1/n^2)$ while each $w_k = 1$. Thus from Theorem 3.1 it follows that $n^{-1}\sum(-1)^{i+j}f(X_i, X_j)$ has a nonnormal limit even if $f$ is nondegenerate.

On the other hand, suppose $a(i,j) = e_i e_j$, where $(n + c\sqrt{n})/2$ of the $e_i$'s are 1 and $(n - c\sqrt{n})/2$ are $-1$. Then $\langle v_i^2(n)\rangle_n \to c^2/n$ and $w_k = 1$, $z_k = c^2$ for all $k > 1$. This gives an example where Theorem 3.1(iii) applies. This case may be compared to Janson (1984) and to the "surface charge" case in Lieb and Lebowitz (1972).

## 4. Incomplete $U$-statistics.

Incomplete $U$-statistics may be analyzed as in Theorems 2.1 and 3.1, the main change being in the normalization constant. Given a bounded weight function $a(i,j)$, define $N^2(n) = \sum_{1 \le i < j \le n} a^2(i,j)$ and let $C(n)$ be the maximum number of nonzero weights in each collection $\{a(i,1),\ldots, a(i,n)\}$, as $i$ ranges from 1 to $n$. Incomplete $U$-statistics are computationally simpler than the full $U$-statistic to the extent that the connectivity $C(n)$ is small compared to $n$. For such a statistic the proper normalizing factor is $N(n)$ rather than $n$. The next theorem assumes that $C(n)$ is $O(n^\alpha)$ for some $\alpha < 1$; when $\alpha = 1$ we are back to Theorem 2.1.

THEOREM 4.1.   *Let $f$, $a$, $X_i$ and $U_n$ be as in Theorem* 2.1. *Suppose there are constants $k$, $K$ and $\alpha$, $0 \le \alpha < 1$, such that $0 < kn^{1+\alpha} \le N^2(n)$ and $C(n) \le Kn^\alpha$. Then $N^{-1}(n)U_n$ is asymptotically normal with variance $I_2$.*

PROOF.   As in Theorem 2.1, we may assume that $f$ is bounded. In the expansion of $E_F((N^{-1}(n)U_n)^m)$, consider all the nonzero terms with graphs having exactly $r$ connected components. (Since $f$ is degenerate, we have

$r \leq m/2$.) There are no more than $n^r C(n)^{m-r}$ of these terms, which is $O(n^{r+\alpha(m-r)})$. On the other hand, $N^{-m}(n)$ is $O(n^{-m(1+\alpha)/2})$. Since $f$ and $a$ are bounded, the sum of all nonzero terms with graphs having exactly $r$ connected components is $O(n^{(r-m/2)(1-\alpha)})$. Thus the only terms which contribute to the limit are those with graphs consisting of two-cycles. The limiting moments can now be found by counting as in Theorem 2.1: All odd moments are zero and the even moments are

$$(4.1) \qquad \mu_{2m} = (2m)! \frac{(2I_2/4)^m}{m!} = I_2^m \frac{(2m)!}{2^m m!}.$$

This completes the proof. $\square$

When $\alpha = 0$, this result is very similar to a theorem proved by Brown and Kildea (1978), which is valid for both degenerate and nondegenerate kernels. Theorem 4.1 is easily extended to nondegenerate kernels with positive $\alpha$ by using the decomposition in (3.5).

**5. Some weighting schemes.** In this section we discuss some weighting functions which are useful in applications and for which the limits $w_k$ can be computed easily.

Perhaps the most common example from physics is the following. Given a bounded sequence of constants $e_1, e_2, \ldots$, define $a(i, j) = e_i e_j$. Since

$$n^{-k} \sum a(i_1, i_2) \cdots a(i_k, i_1) = n^{-k} \sum e_{i_1}^2 \cdots e_{i_k}^2 = \left(\frac{1}{n} \sum e_i^2\right)^k + O\left(\frac{1}{n}\right),$$

we find that $n^{-1}\Sigma e_i^2 \to c > 0$ implies $w_k = c^k$ for all $k > 1$. Thus $n^{-1}U_n/c$ tends to the same limiting distribution as the (unweighted) $U$-statistic in the degenerate kernel case. One has the same limit if the $e_i$ are independent bounded random variables and $n^{-1}\Sigma E(e_i^2) \to c > 0$.

Another useful weight function is one determined by "groups". Let $B = (b_{ij})$ be a real symmetric $s$ by $s$ matrix, and let $g: \mathbf{N} \to \{1, \ldots, s\}$ be a "group identity function." Now we can define a weight function $a$ where the weights are determined by the groups

$$(5.1) \qquad a(i, j) = b_{g(i), g(j)}.$$

For example, if $B$ is diagonal, then the only nonzero terms in $U_n$ are those for which $g(i) = g(j)$ (and thus $U_n$ is the sum of $s$ independent parts.) Suppose now that $r_i(n)$ denotes the proportion of the integers in $\{1, \ldots, n\}$ which are mapped by $g$ into the $i$th group. If $r_i(n) \to r_i$, $1 \leq i \leq s$, then it is easy to show that

$$(5.2) \qquad w_k = \sum_{i_1=1}^{s} \cdots \sum_{i_k=1}^{s} r_{i_1} \cdots r_{i_k} b_{i_1 i_2} b_{i_2 i_3} \cdots b_{i_k i_1}.$$

Thus the cumulants in Theorem 2.1 are easily computed. For example, if $B$ is

diagonal, $w_k = \Sigma r_i^k b_{ii}^k$, and the cumulant generating function is the sum of the cumulant generating functions of the $s$ independent groups.

The statement of Theorem 3.1 may also be modified in this special case. Let $\mathbf{r}(n) = (r_1(n), \ldots, r_s(n))$ and suppose $\mathbf{r}(n) \to \mathbf{r}$ as $n \to \infty$. (It is no restriction to assume that no component of $\mathbf{r}$ is 0.) Then with nondegenerate kernel, asymptotic normality holds when $B\mathbf{r} \neq \mathbf{0}$, and the limit is nonnormal when $nB\mathbf{r}(n)$ is bounded.

The assignment to groups may be made randomly. Suppose $g(i)$ is a sequence of iid random variables that take the value $i$ with probability $r_i$, $1 \le i \le s$. The independence hypothesis of Corollary 2.2 is satisfied, and it can be shown without difficulty that

$$(5.3) \qquad \overline{w}_k = \sum_{i_1=1}^{s} \cdots \sum_{i_k=1}^{s} r_{i_1} \cdots r_{i_k} b_{i_1 i_2} b_{i_2 i_3} \cdots b_{i_k i_1}.$$

Given a symmetric Riemann-integrable function $g: [0,1]^2 \to \mathbf{R}$, consider a triangular scheme of weights, $a_n(i, j) = g(i/n, j/n)$. The proof of Theorem 2.1 is easily adapted to show convergence of the normalized weighted $U$-statistic to a limit with cumulants $(k-1)! w_k I_k / 2$, where

$$w_k = \int_{[0,1]^k} g(x_1, x_2) \cdots g(x_k, x_1) \, dx_1 \cdots dx_k.$$

Theorem 2.1 shows that the limiting distribution is the same as that of the $U$-statistic $n^{-1}\Sigma_{i<j} h(Z_i, Z_j)$ where $Z_i = (X_i, Y_i)$, the $X_i$ are as before and the $Y_i$ are independent random variables uniformly distributed over $[0,1]$, and the kernel $h$ is given by $h((x_1, y_1), (x_2, y_2)) = g(y_1, y_2) f(x_1, x_2)$. This distribution is in turn the same as that of a multiple Wiener integral, following Dynkin and Mandelbaum (1983). It seems a plausible conjecture therefore that the limiting distributions of weighted $U$-statistics of order greater than 2 may be given by corresponding multiple Weiner integrals when the weights have this special form.

## REFERENCES

BLOM, G. (1976). Some properties of incomplete $U$-statistics. *Biometrika* **63** 573–580.

BROWN, B. M. and KILDEA, D. G. (1978). Reduced $U$-statistics and the Hodges–Lehmann estimator. *Ann. Statist.* **6** 828–835.

DYNKIN, E. B. and MANDELBAUM, A. (1983). Symmetric statistics, Poisson point process and multiple Wiener integrals. *Ann. Statist.* **11** 739–745.

FELLER, W. (1968). *An Introduction to Probability Theory and Its Applications*, 3rd ed., **2** 228. Wiley, New York.

GREGORY, G. G. (1977). Large sample theory for $U$-statistics and tests of fit. *Ann. Statist.* **5** 110–123.

HOEFFDING, W. (1948). A class of statistics with asymptotically normal distribution. *Ann. Math. Statist.* **19** 293–325.

JAMMALAMADAKA, S. R. and JANSON, S. (1986). Limit theorems for a triangular scheme of *U*-statistics with applications to inter-point distances. *Ann. Probab.* **14** 1347–1358.

JANSON, S. (1984). The asymptotic distributions of incomplete *U*-statistics. *Z. Wahrsch. Verw. Gebiete* **66** 495–505.

LEE, A. J. (1990). *U-Statistics*. Dekker, New York.

LIEB, E. H. and LEBOWITZ, J. L. (1972). The constitution of matter: Existence of thermodynamics for systems composed of electrons and nuclei. *Adv. Math.* **9** 316–398.

O'NEIL, K. A. and REDNER, R. A. (1991). On the limiting distribution of pair-summable potential functions in many-particle systems. *J. Statist. Phys.* **62** 399–410.

O'NEIL, K. A. and REDNER, R. A. (1992). On the limiting cumulants of *U*-statistics. *Appl. Math. Lett.* **5** 37–40.

RUBIN, H. and VITALE, R. A. (1980). Asymptotic distribution of symmetric statistics. *Ann. Statist.* **8** 165–170.

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

SHAPIRO, C. P. and HUBERT, L. (1979). Asymptotic normality of permutation statistics derived from weighted sums of bivariate functions. *Ann. Statist.* **7** 788–794.

DEPARTMENT OF MATHEMATICAL
AND COMPUTER SCIENCES
UNIVERSITY OF TULSA
600 S. COLLEGE AVENUE
TULSA, OKLAHOMA 74104