

A COMBINATORIAL CENTRAL LIMIT THEOREM¹

BY WASSILY Hoeffding

Institute of Statistics, University of North Carolina

1. Summary. Let (Y_{n1}, \dots, Y_{nn}) be a random vector which takes on the $n!$ permutations of $(1, \dots, n)$ with equal probabilities. Let $c_n(i, j)$, $i, j = 1, \dots, n$, be n^2 real numbers. Sufficient conditions for the asymptotic normality of

$$S_n = \sum_{i=1}^n c_n(i, Y_{ni})$$

are given (Theorem 3). For the special case $c_n(i, j) = a_n(i)b_n(j)$ a stronger version of a theorem of Wald, Wolfowitz and Noether is obtained (Theorem 4). A condition of Noether is simplified (Theorem 1).

2. Introduction and statement of results. An example of what is here meant by a combinatorial central limit theorem is a solution of the following problem. For every positive integer n there are given $2n$ real numbers $a_n(i)$, $b_n(i)$, $i = 1, \dots, n$. It is assumed that the $a_n(i)$ are not all equal and the $b_n(i)$ are not all equal. Let (Y_{n1}, \dots, Y_{nn}) be a random vector which takes on the $n!$ permutations of $(1, \dots, n)$ with equal probabilities $1/n!$. Under what conditions is

$$(1) \quad S_n = \sum_{i=1}^n a_n(i)b_n(Y_{ni})$$

asymptotically normally distributed as $n \rightarrow \infty$?

Throughout this paper a random variable S_n will be called asymptotically normal or asymptotically normally distributed if

$$\lim_{n \rightarrow \infty} \Pr\{S_n - ES_n \leq x \sqrt{\text{var} S_n}\} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-1/2 y^2} dy, \quad -\infty < x < \infty,$$

where ES_n and $\text{var} S_n$ are the mean and the variance of S_n .

In the particular case $a_n(i) = b_n(i) = i$ the asymptotic normality of S_n was proved by Hotelling and Pabst [2]. The first general result is due to Wald and Wolfowitz [6], who showed that S_n is asymptotically normal if, as $n \rightarrow \infty$,

$$(2) \quad \frac{\frac{1}{n} \sum_{i=1}^n (a_n(i) - \bar{a}_n)^r}{\left[\frac{1}{n} \sum_{i=1}^n (a_n(i) - \bar{a}_n)^2 \right]^{r/2}} = O(1), \quad r = 3, 4, \dots,$$

and

$$(3) \quad \frac{\frac{1}{n} \sum_{i=1}^n (b_n(i) - \bar{b}_n)^r}{\left[\frac{1}{n} \sum_{i=1}^n (b_n(i) - \bar{b}_n)^2 \right]^{r/2}} = O(1), \quad r = 3, 4, \dots,$$

¹ Work done under the sponsorship of the Office of Naval Research.

where

$$\bar{a}_n = \frac{1}{n} \sum_{i=1}^n a_n(i), \quad \bar{b}_n = \frac{1}{n} \sum_{i=1}^n b_n(i).$$

Noether [5] proved that condition (3) can be replaced by the weaker condition

$$(4) \quad \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n (b_n(i) - \bar{b}_n)^r}{\left[\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2 \right]^{r/2}} = 0, \quad r = 3, 4, \dots$$

This condition can be simplified as follows.

THEOREM 1. Condition (4) is equivalent to either of the following two conditions:

$$(5) \quad \lim_{n \rightarrow \infty} \frac{\sum_{i=1}^n |b_n(i) - \bar{b}_n|^r}{\left[\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2 \right]^{r/2}} = 0 \quad \text{for some } r > 2;$$

$$(6) \quad \lim_{n \rightarrow \infty} \frac{\max_{1 \leq i \leq n} (b_n(i) - \bar{b}_n)^2}{\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2} = 0.$$

Hence conditions (2) and (5) or (2) and (6) are sufficient for the asymptotic normality of (1).

The proof is given in Section 3. For a more general condition and a stronger but simpler condition see Theorem 4 below.

One extension of this problem was considered by Daniels [1], who studied the asymptotic distribution of

$$\sum_{i=1}^n \sum_{j=1}^n a_n(i, j) b_n(Y_{ni}, Y_{nj}).$$

The present paper is concerned with an alternative extension. It considers the distribution of

$$(7) \quad S_n = \sum_{i=1}^n c_n(i, Y_{ni}),$$

where $c_n(i, j)$, $i, j = 1, \dots, n$, are n^2 real numbers, defined for every positive integer n . In the particular case $c_n(i, j) = a_n(i)b_n(j)$, (7) reduces to (1).

Let

$$(8) \quad d_n(i, j) = c_n(i, j) - \frac{1}{n} \sum_{g=1}^n c_n(g, j) - \frac{1}{n} \sum_{h=1}^n c_n(i, h) + \frac{1}{n^2} \sum_{g=1}^n \sum_{h=1}^n c_n(g, h).$$

THEOREM 2. The mean and variance of

$$S_n = \sum_{i=1}^n c_n(i, Y_{ni})$$

are

$$(9) \quad ES_n = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n c_n(i, j),$$

$$(10) \quad \text{var } S_n = \frac{1}{n-1} \sum_{i=1}^n \sum_{j=1}^n d_n^2(i, j).$$

Henceforth we assume that $d_n(i, j) \neq 0$ for some (i, j) , so that $\text{var } S_n > 0$. In the special case $c_n(i, j) = a_n(i)b_n(j)$ this corresponds to the assumption that the $a_n(i)$ are not all equal and the $b_n(j)$ are not all equal.

THEOREM 3. *The distribution of $S_n = \sum_{i=1}^n c_n(i, Y_{ni})$ is asymptotically normal if*

$$(11) \quad \lim_{n \rightarrow \infty} \frac{\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d_n^r(i, j)}{\left[\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d_n^2(i, j) \right]^{r/2}} = 0, \quad r = 3, 4, \dots$$

Condition (11) is satisfied if

$$(12) \quad \lim_{n \rightarrow \infty} \frac{\max_{1 \leq i, j \leq n} d_n^2(i, j)}{\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d_n^2(i, j)} = 0.$$

Theorems 2 and 3 will be proved in Sections 4 and 5.

For the special case $c_n(i, j) = a_n(i)b_n(j)$, Theorem 3 immediately gives

THEOREM 4. *The distribution of $S_n = \sum_{i=1}^n a_n(i)b_n(Y_{ni})$ is asymptotically normal if*

$$(13) \quad \lim_{n \rightarrow \infty} n^{\frac{1}{2}r-1} \frac{\sum_{i=1}^n (a_n(i) - \bar{a}_n)^r \sum_{i=1}^n (b_n(i) - \bar{b}_n)^r}{\left[\sum_{i=1}^n (a_n(i) - \bar{a}_n)^2 \right]^{r/2} \left[\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2 \right]^{r/2}} = 0, \quad r = 3, 4, \dots$$

Condition (13) is satisfied if

$$(14) \quad \lim_{n \rightarrow \infty} n \frac{\max_{1 \leq i \leq n} (a_n(i) - \bar{a}_n)^2}{\sum_{i=1}^n (a_n(i) - \bar{a}_n)^2} \frac{\max_{1 \leq i \leq n} (b_n(i) - \bar{b}_n)^2}{\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2} = 0.$$

It will be observed that the symmetrical condition (13) contains Noether's condition (2) and (4) as a special case.

Let $X_n = (X_{n1}, \dots, X_{nn})$ be independent of and have the same distribution as $Y_n = (Y_{n1}, \dots, Y_{nn})$.

THEOREM 5. *The random variable*

$$(15) \quad S'_n = \sum_{i=1}^n c_n(X_{ni}, Y_{ni})$$

has the same distribution as S_n in (7).

In fact, the conditional distribution of S'_n given that $X_n = p$, a fixed permutation of $(1, \dots, n)$, is independent of p because the distribution of Y_n is invariant under permutations of its components.

The distribution of sums of the form (1) has attracted the attention of statisticians in connection with nonparametric tests (see, for example, [2], [6], [3]) and sampling from a finite population (which leads to the case $a_n(i) = 0$ for $i > m$; cf. also Madow [4]). More general sums of the form (7) or (15) are likewise of interest in nonparametric theory. Thus it follows from results of Lehmann and Stein [3] that a test of the hypothesis that U_1, \dots, U_n are independent and identically distributed, which is most powerful similar against the alternative that the joint frequency function is $f_1(u_1) \cdots f_n(u_n)$ is based on a statistic of the form (7) with

$$c_n(i, j) = \log f_i(u_j),$$

where the u_j are the observed sample values. If the n pairs $(U_1, V_1), \dots, (U_n, V_n)$ are independent and identically distributed, a test of the hypothesis that U_i and V_i are independent which is most powerful similar against the alternative that their joint frequency function is $f(u, v)$ is based on a statistic of the form (15) with $c_n(i, j) = \log f(u_i, v_j)$, where $(u_1, v_1), \dots, (u_n, v_n)$ are the observed values.

In these examples the numbers $c_n(i, j)$ are random variables. An application of some of the present results to such cases will be considered by the author in a forthcoming paper.

3. Proof of Theorem 1. Let

$$g_i = \frac{b_n(i) - \bar{b}_n}{\left[\sum_{i=1}^n (b_n(i) - \bar{b}_n)^2 \right]^{1/2}},$$

$$G_n = \max(|g_1|, \dots, |g_n|).$$

Theorem 1 asserts the equivalence of the three relations

(16) $\lim_{n \rightarrow \infty} \sum_{i=1}^n g_i^r = 0, \quad r = 3, 4, \dots;$

(17) $\lim_{n \rightarrow \infty} \sum_{i=1}^n |g_i|^r = 0 \quad \text{for some } r > 2;$

(18) $\lim_{n \rightarrow \infty} G_n = 0.$

We have

$$\sum_{i=1}^n g_i^2 = 1,$$

and hence for $r > 2$

$$G_n^r \leq \sum_{i=1}^n |g_i|^r \leq G_n^{r-2} \sum_{i=1}^n g_i^2 = G_n^{r-2}.$$

The equivalence of (16), (17) and (18) follows immediately.

4. Proof of Theorem 2. The subscript n in Y_{ni} , $c_n(i, j)$, etc., will henceforth be omitted. We note that if the subscripts i_1, \dots, i_m are distinct, the expected value of a function $f(Y_{i_1}, \dots, Y_{i_m})$ is equal to

$$\frac{1}{n(n-1)\dots(n-m+1)} \sum'_{j_1, \dots, j_m} f(j_1, \dots, j_m),$$

where the sum Σ' is extended over all m -tuples (j_1, \dots, j_m) of distinct integers from 1 to n . Relation (9) follows immediately.

Let

$$(19) \quad T_n = \sum_{i=1}^n d(i, Y_i),$$

where $d(i, j) = d_n(i, j)$ is defined by (8). Using (9), we get

$$(20) \quad T_n = S_n - ES_n.$$

Also

$$(21) \quad \sum_{i=1}^n d(i, j) = 0 \text{ for all } j, \quad \sum_{j=1}^n d(i, j) = 0 \text{ for all } i.$$

Hence

$$Ed(i, Y_i) = 0,$$

$$Ed^2(i, Y_i) = \frac{1}{n} \sum_{j=1}^n d^2(i, j),$$

and if $i \neq j$,

$$Ed(i, Y_i)d(j, Y_j) = \frac{1}{n(n-1)} \sum'_{g,h} d(i, g)d(j, h)$$

$$= \frac{-1}{n(n-1)} \sum_{g=1}^n d(i, g)d(j, g).$$

Therefore

$$\text{var } S_n = \text{var } T_n = \sum_{i=1}^n Ed^2(i, Y_i) + \sum'_{i,j} Ed(i, Y_i)d(j, Y_j)$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d^2(i, j) - \frac{1}{n(n-1)} \sum_{g=1}^n \sum'_{i,j} d(i, g)d(j, g)$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d^2(i, j) + \frac{1}{n(n-1)} \sum_{g=1}^n \sum_{i=1}^n d^2(i, g),$$

which gives relation (10).

5. Proof of Theorem 3. Let

$$(22) \quad M_{r,n} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d^r(i, j),$$

$$(23) \quad \bar{M}_{r,n} = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n |d(i, j)|^r,$$

$$(24) \quad D_n = \max_{1 \leq i, j \leq n} |d(i, j)|.$$

Then $\text{var } S_n = n/(n-1) M_{2,n}$. Since, by hypothesis, $\text{var } S_n > 0$, we may and shall assume that

$$(25) \quad M_{2,n} = 1.$$

Conditions (11) and (12) can now be written as

$$(26) \quad \lim_{n \rightarrow \infty} M_{r,n} = 0, \quad r = 3, 4, \dots,$$

and

$$(27) \quad \lim_{n \rightarrow \infty} D_n = 0.$$

That (27) implies (26) is seen from the inequalities

$$|M_{r,n}| \leq \bar{M}_{r,n} \leq D_n^{r-2} M_{2,n} = D_n^{r-2} \quad \text{for } r > 2.$$

Since

$$\bar{M}_{2k+1,n}^2 \leq M_{2k,n} M_{2k+2,n}, \quad k = 1, 2, \dots,$$

condition (26) implies

$$(28) \quad \lim_{n \rightarrow \infty} \bar{M}_{r,n} = 0, \quad r = 3, 4, \dots$$

As $\text{var } S_n \rightarrow 1$, it is now sufficient to demonstrate that under conditions (25) and (28), $T_n = S_n - ES_n$ has a normal limiting distribution with mean 0 and variance 1. This will be proved by showing that

$$(29) \quad \lim_{n \rightarrow \infty} ET_n^r = \begin{cases} 1 \cdot 3 \cdots (r-1) & \text{if } r \text{ is even,} \\ 0 & \text{if } r \text{ is odd.} \end{cases}$$

The r th moment of T_n ,

$$(30) \quad ET_n^r = E \sum_{i_1=1}^n \cdots \sum_{i_r=1}^n d(i_1, Y_{i_1}) \cdots d(i_r, Y_{i_r}),$$

can be written as a sum of terms of the form

$$(31) \quad I(r, e_1, \dots, e_m) = \sum'_{i_1, \dots, i_m} E d^{e_1}(i_1, Y_{i_1}) \cdots d^{e_m}(i_m, Y_{i_m}),$$

where $e_i \geq 1, e_1 + \dots + e_m = r$. The number of terms (31) is independent of n . It will be shown that

$$(32) \quad \lim_{n \rightarrow \infty} I(r, e_1, \dots, e_m) = 0 \quad \text{unless } r = 2m, \quad e_1 = \dots = e_m = 2,$$

$$(33) \quad \lim_{n \rightarrow \infty} I(r, 2, \dots, 2) = 1 \quad \text{if } r \text{ even,}$$

and that the number of terms $I(r, 2, \dots, 2)$ in (30) with r even equals $1 \cdot 3 \cdots (r - 1)$. Then (29) holds, and the theorem will be proved.

We have for $n \rightarrow \infty$

$$(34) \quad I(r, e_1, \dots, e_m) \sim n^{-m} \sum'_{i_1, \dots, i_m} \sum'_{j_1, \dots, j_m} d^{e_1}(i_1, j_1) \cdots d^{e_m}(i_m, j_m).$$

The right-hand side can be written as a sum of terms which, apart from the sign, are of the form

$$(35) \quad n^{-m} J(r, p, q, e_1, \dots, e_m) = n^{-m} \sum_{i_1=1}^n \cdots \sum_{i_p=1}^n \sum_{j_1=1}^n \cdots \sum_{j_q=1}^n d^{e_1}(i_{c_1}, j_{d_1}) \cdots d^{e_m}(i_{c_m}, j_{d_m}),$$

where

$$1 \leq p \leq m, \quad 1 \leq q \leq m, \\ 1 \leq c_g \leq p, \quad 1 \leq d_h \leq q, \quad (g, h = 1, \dots, m),$$

and for every integer u , $1 \leq u \leq p$ ($1 \leq u \leq q$) at least one c_g (d_h) is equal to u . The number of terms (35) is independent of n .

The sum J in (35) can be written as a product of $s \geq 1$ sums of a similar form,

$$(36) \quad J(r, p, q, e_1, \dots, e_m) = \prod_{k=1}^s J(r_k, p_k, q_k, e_{k1}, \dots, e_{km_k}),$$

where

$$(e_{k1}, \dots, e_{km_k}), \quad k = 1, \dots, s,$$

are s disjoint subsets of (e_1, \dots, e_m) ,

$$(37) \quad \begin{aligned} e_{k1} + \cdots + e_{km_k} &= r_k, & r_1 + \cdots + r_s &= r, \\ p_1 + \cdots + p_s &= p, & q_1 + \cdots + q_s &= q, \\ m_1 + \cdots + m_s &= m. \end{aligned}$$

We observe that

$$(38) \quad 1 \leq p_k \leq m_k, \quad 1 \leq q_k \leq m_k, \quad m_k \leq r_k.$$

It will be assumed that s is the greatest possible number of factors into which $J(r, p, q, e_1, \dots, e_m)$ can be decomposed in the form (36). If $s = 1$, the number of equalities between the subscripts c or between the subscripts d in (35) must be at least $m - 1$. The total number of subscripts c, d being $2m$, there are at most $m + 1$ distinct subscripts, so that $p + q \leq m + 1$. If

$$(39) \quad (c_g, d_g) = (c_h, d_h) \quad \text{for some } (g, h), g \neq h,$$

we have strict inequality. For an arbitrary s we have in a similar way

$$(40) \quad p_k + q_k \leq m_k + 1, \quad k = 1, \dots, s,$$

and hence

$$(41) \quad p + q \leq m + s,$$

with strict inequality in the case (39).

By Hölder's inequality, from (35),

$$\begin{aligned} |J(r, p, q, e_1, \dots, e_m)| &\leq \prod_{g=1}^m \left(\sum_{i_1} \dots \sum_{i_p} \sum_{j_1} \dots \sum_q |d(i_{c_g}, j_{d_g})|^r \right)^{e_g/r} \\ &= \prod_{g=1}^m (n^{p+q-1} \bar{M}_{r,n})^{e_g/r} = n^{p+q-1} \bar{M}_{r,n}. \end{aligned}$$

Similarly,

$$|J(r_k, p_k, q_k, e_{k1}, \dots, e_{km_k})| \leq n^{p_k+q_k-1} \bar{M}_{r_k,n}.$$

Hence, by (36),

$$(42) \quad n^{-m} |J(r, p, q, e_1, \dots, e_m)| \leq n^{p+q-s-m} \bar{M}_{r_1,n} \dots \bar{M}_{r_s,n}.$$

If, for some k , $r_k = 1$, then, by (38) and (37), $p_k = q_k = m_k = e_{k1} = 1$, and hence $J = 0$ by (21). Thus we may assume $r_k \geq 2$, $k = 1, \dots, s$. Then, by (28), $\bar{M}_{r_1,n} \dots \bar{M}_{r_s,n} \rightarrow 0$ unless $r_1 = \dots = r_s = 2$. It now follows from (42) and (41) that

$$(43) \quad \lim_{n \rightarrow \infty} n^{-m} J(r, p, q, e_1, \dots, e_m) = 0$$

except perhaps when $r_1 = \dots = r_s = 2$.

If $r_1 = \dots = r_s = 2$, we have

$$(44) \quad n^{-m} J(r, p, q, e_1, \dots, e_m) = O(n^{p+q-s-m}).$$

By (38), $r_k = 2$ implies $m_k = 1$ or 2 . If $m_k = 2$, then $e_{k1} = e_{k2} = 1$ and $p_k + q_k \leq 3$ by (40). If $p_k + q_k = 3$, the corresponding J -factor is of the form

$$\sum_i \sum_j \sum_k d(i, j)d(i, k) \quad \text{or} \quad \sum_i \sum_j \sum_k d(i, k)d(j, k),$$

both of which vanish by (21). If $m_k = 2$ and $p_k + q_k = 2$, we have case (39) and hence, by the remark following (41), $p + q - s - m < 0$. By (44), this implies (43).

Thus the only case where (43) need not hold is $r_k = 2, m_k = 1$ for $k = 1, \dots, s$. Then $p_k = q_k = 1, e_{k1} = 2$, hence

$$r = 2s = 2m, \quad p = q = r/2$$

$$e_1 = \dots = e_m = 2.$$

This proves relation (32), and (33) follows from

$$\begin{aligned} I(r, 2, \dots, 2) &\sim n^{-r/2} J\left(r, \frac{r}{2}, \frac{r}{2}, 2, \dots, 2\right) \\ &= n^{-r/2} [J(2, 1, 1, 2)]^{r/2} \\ &= M_{2,n}^{r/2} = 1. \end{aligned}$$

It remains to determine the number of terms $I(r, 2, \dots, 2)$ in (30) when r is even. This is the number of ways the subscripts i_1, \dots, i_r can be tied in $r/2$ groups of two, which is $(r - 1)(r - 3) \cdots 3 \cdot 1$. The proof is complete.

REFERENCES

- [1] H. E. DANIELS, "The relation between measures of correlation in the universe of sample permutations," *Biometrika*, Vol. 33 (1944), pp. 129-135.
- [2] H. HOTELLING AND M. PABST, "Rank correlation and tests of significance involving no assumption of normality," *Annals of Math. Stat.*, Vol. 7 (1936), pp. 29-43.
- [3] E. L. LEHMANN AND C. STEIN, "On the theory of some nonparametric hypotheses," *Annals of Math. Stat.*, Vol. 20 (1949), pp. 28-45.
- [4] W. G. MADOW, "On the limiting distributions of estimates based on samples from finite universes," *Annals of Math. Stat.*, Vol. 19 (1948), pp. 535-545.
- [5] G. E. NOETHER, "On a theorem by Wald and Wolfowitz," *Annals of Math. Stat.*, Vol. 20 (1949), pp. 455-458.
- [6] A. WALD AND J. WOLFOWITZ, "Statistical tests based on permutations of the observations," *Annals of Math. Stat.*, Vol. 15 (1944), pp. 358-372.