

UNIVERSITY OF NORTH CAROLINA
Department of Statistics
Chapel Hill, N. C.

Mathematical Sciences Directorate
Air Force Office of Scientific Research
Washington 25, D. C.

AFOSR Report No.

THE STRONG LAW OF LARGE NUMBERS FOR U-STATISTICS

by

Wassily Hoeffding

July, 1961

Contract No. AF 49(638)-261

Let X_1, X_2, \dots be independent, identically distributed random variables, let $f(X_1, \dots, X_r)$ be a function whose expected value, θ , exists, and let \bar{F}_n be the arithmetic (with equal weights) of the values $f(X_{i_1}, \dots, X_{i_r})$, for all r -tuples (i_1, \dots, i_r) of distinct positive integers not exceeding n . It is shown that $\bar{F}_n \longrightarrow \theta$ almost surely.

Qualified requestors may obtain copies of this report from the ASTIA Document Service Center, Arlington Hall Station, Arlington 12, Virginia. Department of Defense contractors must be established for ASTIA services, or have their "need-to-know" certified by the cognizant military agency of their project or contract.

Institute of Statistics
Mimeograph Series No. 302

THE STRONG LAW OF LARGE NUMBERS FOR U-STATISTICS

by Wassily Hoeffding¹

University of North Carolina

1. Introduction. Let X_1, X_2, \dots be a sequence of mutually independent, identically distributed random variables taking values in a space \mathcal{X} . Let r be a positive integer and f a real-valued measurable function on \mathcal{X}^r . For $n \geq r$ define \bar{f}_n as the arithmetic mean

$$(1) \quad \bar{f}_n = \frac{1}{n(n-1)\dots(n-r+1)} \sum f(X_{i_1}, \dots, X_{i_r}),$$

where the sum is extended over all r -tuples i_1, \dots, i_r of distinct positive integers not exceeding n . The following theorem will be proved. (We write a.s. for almost surely or almost sure.)

Theorem. If $E [|f(X_1, \dots, X_r)|] < \infty$, then

$$(2) \quad \bar{f}_n \longrightarrow E [f(x_1, \dots, x_r)] \quad \text{a.s.}$$

The theorem contains, for $r = 1$, the sufficiency part of Kolmogorov's strong law of large numbers. The asymptotic behavior of random variables of the form (1) (which have been called U-statistics) has been studied, among others, by the author in [2], where a central limit theorem for such sums has been proved, and by P. K. Sen [3], who obtained the result (2) under the assumption that a moment of order higher than $2 - r^{-1}$ is finite.

The proof of the theorem makes use of the fact that \bar{f}_n can be represented as a linear combination of r random variables each of which has a martingale

1

This research was supported by the United States Air Force through the Air Force Office of Scientific Research of the Air Research and Development Command, under Contract No. AF 49(638)-261. Reproduction in whole or in part is permitted for any purpose of the United States Government.

property (Lemma 1). The summands are suitably truncated and a version of Doob's extension to martingales of Kolmogorov's inequality (see Lemma 2) is applied.

2. Preliminaries. We first observe that since the sum in (1) is symmetric, each term may be replaced by the arithmetic mean of the $r!$ terms whose subscripts are permutations of the same set of integers. We thus may, and shall assume that f is invariant under permutations:

$$(3) \quad f(x_{i_1}, \dots, x_{i_r}) = f(x_1, \dots, x_r)$$

for all permutations (i_1, \dots, i_r) of $(1, \dots, r)$. We now can write

$$(4) \quad \bar{F}_n = \binom{n}{r}^{-1} S_n, \quad S_n = \sum_{1 \leq i_1 < \dots < i_r \leq n} f(X_{i_1}, \dots, X_{i_r})$$

We shall also assume with no loss of generality that $E [f(X_1, \dots, X_r)] = 0$ or

$$(5) \quad \int \dots \int f(x_1, \dots, x_r) dF(x_1) \dots dF(x_r) = 0,$$

where F denotes the common probability measure of the random variables X_n .

Lemma 1. If condition (5) is satisfied, we have

$$(6) \quad \binom{n}{r}^{-1} S_n = \sum_{h=1}^r \binom{r}{h} \binom{n}{h}^{-1} S_{h,n},$$

where

$$(7) \quad S_{h,n} = \sum_{1 \leq i_1 < \dots < i_h \leq n} f_h(X_{i_1}, \dots, X_{i_h}),$$

the function f_h ($h = 1, \dots, r$) is invariant under permutations of its arguments and satisfies the condition

$$(8) \quad \int f_h(x_1, \dots, x_{h-1}, x) dF(x) \equiv 0;$$

and $S_{h,n}$ has the martingale property

$$(9) \quad E [S_{h,n} | X_1, \dots, X_m] = S_{h,m}, \quad h \leq m \leq n.$$

Proof. We shall write $f(x_1, \dots, x_{r-1}, *)$ for $\int f(x_1, \dots, x_{r-1}, x) dF(x)$, etc. We define the functions f_h as follows.

$$\begin{aligned} f_1(x_1) &= f(x_1, *, \dots, *) \\ f_2(x_1, x_2) &= f(x_1, x_2, *, \dots, *) - f_1(x_1) - f_1(x_2) \\ (10) \quad f_3(x_1, x_2, x_3) &= f(x_1, x_2, x_3, *, *, *) - f_1(x_1) - f_1(x_2) - f_1(x_3) \\ &\quad - f_2(x_1, x_2) - f_2(x_1, x_3) - f_2(x_2, x_3) \\ &\quad \dots \dots \dots \\ f_r(x_1, \dots, x_r) &= f(x_1, \dots, x_r) - f_1(x_1) - \dots - f_1(x_r) \\ &\quad - f_2(x_1, x_2) - f_2(x_1, x_3) - \dots - f_2(x_{r-1}, x_r) \\ &\quad \dots - f_{r-1}(x_1, \dots, x_{r-1}) - \dots - f_{r-1}(x_2, \dots, x_r). \end{aligned}$$

The last equation expresses f in terms of f_1, \dots, f_r . Inserting this expression in (4), we obtain equation (6). The invariance of f_h under permutations and equation (8) follow from the definition of the f_h and assumption (5). Finally, for $h \leq m < n$, $S_{h,n} - S_{h,m}$ is a sum of terms $f_h(X_{i_1}, \dots, X_{i_h})$ with $i_j > m$ for some j . Hence, due to (8), the conditional expectation of $S_{h,n} - S_{h,m}$ when X_1, \dots, X_m are fixed is 0, and the martingale property (9) follows.

Due to Lemma 1 we may assume that

$$(11) \quad \int f(x_1, \dots, x_{r-1}, x) dF(x) \equiv 0$$

and hence

$$(12) \quad E [S_n | X_1, \dots, X_m] = S_m, \quad r \leq m \leq n.$$

The theorem will be proved if we show that $n^{-r} S_n \rightarrow 0$ a.s. To this end we first truncate the random variables $f(x_{i_1}, \dots, x_{i_r})$ as follows. For $j > 0$ let

$$(13) \quad f^{(j)}(x_1, \dots, x_r) = \begin{cases} f(x_1, \dots, x_r) & \text{if } |f(x_1, \dots, x_r)| \leq j^r \\ 0 & \text{otherwise} \end{cases},$$

$$(14) \quad S'_n = \sum_{1 \leq i_1 < \dots < i_r \leq n} f^{(i_r)}(X_{i_1}, \dots, X_{i_r}).$$

It will be shown that $n^{-r}(S_n - S'_n) \rightarrow 0$ a.s. Now S'_n does not have the martingale property. However, by using a device similar to that of Lemma 1, S'_n can be written as $T_n + R_n$, where T_n has the martingale property and R_n is a sum of terms having less than r arguments. By induction on r it is shown that $n^{-r} R_n \rightarrow 0$ a.s., and Lemma 2 below (which, in essence, is known) will imply that $n^{-r} T_n \rightarrow 0$ a.s. The proof will be given only for $r = 2$, in such a way that the extension to the case $r \sim 2$ can be left to the reader.

Lemma 2. Let X_1, X_2, \dots be a sequence of random variables, and for each positive integer n let Z_n be a real-valued function of X_1, \dots, X_n such that, for some $a \geq 1$, $E [|Z_n|^a] < \infty$ and

$$(15) \quad E [Z_n | X_1, \dots, X_m] = Z_m, \quad 1 \leq m \leq n.$$

If, for some $b > 0$,

$$(16) \quad \sum_{k=1}^{\infty} 2^{-abk} E [|Z_k|^a] < \infty ,$$

then $n^{-b} Z_n \longrightarrow 0$ a.s.

Proof. We have for $t > 0$

$$(17) \quad P [|Z_m| \geq t \text{ for some } m \leq n] \leq t^{-a} E [|Z_n|^a] .$$

This is a trivial extension of Doob's generalization ([1], p. 314) of Kolmogorov's inequality. By a standard argument, (16) and (17) imply $n^{-b} Z_n \longrightarrow 0$ a.s.

3. Proof of the theorem for $r = 2$. We now let

$$(18) \quad S_n = \sum_{1 \leq i < j \leq n} f(X_i, X_j) ,$$

where

$$(19) \quad f(x, y) = f(y, x) , \quad \int f(x, y) dF(y) \equiv 0 ,$$

and define

$$(20) \quad f^{(j)}(x, y) = \begin{cases} f(x, y) & , \quad |f(x, y)| \leq j^2 \\ 0 & \text{otherwise} \end{cases} ,$$

$$(21) \quad S'_n = \sum_{1 \leq i < j \leq n} f^{(j)}(X_i, X_j) .$$

Lemma 3. $n^{-2}(S_n - S'_n) \longrightarrow 0$ a.s.

Proof. Let $Y_{ij} = f(X_i, X_j)$, $Y'_{ij} = f^{(j)}(X_i, X_j)$. We have for $i \neq j$

$$\begin{aligned}
 P [Y_{ij} \neq Y'_{ij}] &= \iint_{|f(x,y)| > j^2} dF(x)dF(y) \\
 &= \sum_{v=j^2}^{\infty} \iint_{v < |f(x,y)| \leq v+1} dF(x)dF(y) \leq \sum_{v=j^2}^{\infty} v^{-1} a_{v+1} ,
 \end{aligned}$$

where

$$(22) \quad a_n = \iint_{n-1 < |f(x,y)| \leq n} |f(x,y)| dF(x)dF(y) .$$

Hence

$$\begin{aligned}
 \sum_{1 \leq i < j < \infty} P [Y_{ij} \neq Y'_{ij}] &\leq \sum_{1 \leq i < j < \infty} \sum_{v=j^2}^{\infty} v^{-1} a_{v+1} \\
 &= \sum_{v=4}^{\infty} \sum_{2 \leq j \leq v^{1/2}} (j-1)v^{-1} a_{v+1} \leq \sum_{v=4}^{\infty} a_{v+1} < \infty .
 \end{aligned}$$

This implies

$$(23) \quad \lim_{m \rightarrow \infty} P [Y_{ij} \neq Y'_{ij} \text{ for some } (i,j), m < i < j < \infty] = 0 .$$

Now let m be a fixed integer. For $n > m$,

$$\begin{aligned}
 (24) \quad n^{-2}(S_n - S'_n) &= n^{-2}(S_m - S'_m) + \sum_{i=1}^m n^{-2} \sum_{j=m+1}^n (Y_{ij} - Y'_{ij}) \\
 &\quad + n^{-2} \sum_{m < i < j \leq n} (Y_{ij} - Y'_{ij}) .
 \end{aligned}$$

By (23) we can choose m so large that with a probability arbitrarily close to 1 the last sum in (24) is 0 for all n . Clearly $n^{-2}(S_m - S'_m) \rightarrow 0$ a.s. as $n \rightarrow \infty$. It remains to show that for a fixed $i \leq m$, $n^{-2} \sum_{j=m+1}^n (Y_{ij} - Y'_{ij}) \rightarrow 0$ a.s. It will suffice to prove this for $i = m = 1$.

Let $g^{(j)}(x,y) = f(x,y)$ if $|f(x,y)| > j^2$, $g^{(j)}(x,y) = 0$ otherwise. Then

$$(25) \quad n^{-2} \sum_{j=2}^n (Y_{1j} - Y'_{1j}) = n^{-2} \sum_{j=2}^n g^{(j)}(X_1, X_j) = n^{-2} V_n + n^{-2} \sum_{j=2}^n g^{(j)}(X_1, *) ,$$

where $V_n = \sum_{j=2}^n [g^{(j)}(X_1, X_j) - g^{(j)}(X_1, *)]$. Since $|g^{(j)}(x, *)| \leq \int |f(x,y)| dF(y)$, it is clear that the last term in (25) tends to 0 a.s. Also $E[V_n | X_1, \dots, X_m] = V_m$ for $1 < m < n$. We apply Lemma 2 with $Z_n = V_n$, $a = 1$, $b = 2$. Since $E[|V_n|] \leq c n$, condition (16) is satisfied and hence $n^{-2} V_n \rightarrow 0$ a.s. The proof of Lemma 3 is complete.

We now show that $n^{-2} S'_n \rightarrow 0$ a.s. Let

$$(26) \quad f_2^{(j)}(x,y) = f^{(j)}(x,y) - f^{(j)}(x,*) - f^{(j)}(y,*) + f^{(j)}(*,*) .$$

Then

$$(27) \quad S'_n = \sum_{1 \leq i < j \leq n} f_2^{(j)}(X_i, X_j) + \sum_{1 \leq i < j \leq n} [f^{(j)}(X_i, *) + f^{(j)}(X_j, *)] - \sum_{j=2}^n (j-1) f^{(j)}(*,*) .$$

Now

$$\begin{aligned} |f^{(j)}(*,*)| &= \left| \iint_{|f(x,y)| \leq j^2} f(x,y) dF(x) dF(y) \right| \\ &= \left| \iint_{|f(x,y)| > j^2} f(x,y) dF(x) dF(y) \right| \leq \sum_{v=j^2+1}^{\infty} a_v . \end{aligned}$$

Hence

$$\begin{aligned}
\left| \sum_{j=2}^n (j-1) f^{(j)}(*,*) \right| &\leq \sum_{j=2}^n (j-1) \sum_{v=j^2+1}^{\infty} a_v \\
&\leq \sum_{v=5}^{\infty} \sum_{2 \leq j \leq \min(n, \sqrt{v})} (j-1) a_v \leq \sum_{v=5}^n n a_v + \sum_{v=n+1}^{\infty} n^2 a_v,
\end{aligned}$$

so that

$$(28) \quad n^{-2} \sum_{j=2}^n (j-1) f^{(j)}(*,*) \longrightarrow 0.$$

Next, making use of (19),

$$f^{(j)}(x,*) = \int_{|f(x,y)| \leq j^2} f(x,y) dF(y) - \int_{|f(x,y)| > j^2} f(x,y) dF(y),$$

so that

$$(29) \quad |f^{(j)}(x,*)| \leq \int_{|f(x,y)| > j^2} |f(x,y)| dF(y) = h_j(x), \quad \text{say.}$$

We note that $h_j(x) \geq 0$, $h_j(x) \leq h_m(x)$ for $j > m$, and $\int h_m(x) dF(x) \longrightarrow 0$ as $m \longrightarrow \infty$. Now for $n > m$,

$$\begin{aligned}
\left| \sum_{1 \leq i < j \leq n} [f^{(j)}(X_i,*) + f^{(j)}(X_j,*)] \right| &\leq \sum_{1 \leq i < j \leq m} [h_0(X_i) + h_0(X_j)] \\
&\quad + \sum_{i=1}^m \sum_{j=m+1}^n [h_0(X_i) + h_m(X_j)] + \sum_{m < i < j \leq n} [h_m(X_i) + h_m(X_j)] \\
&= (n-1) \sum_{i=1}^m h_0(X_i) + (n-1) \sum_{i=m+1}^n h_m(X_i).
\end{aligned}$$

For m fixed and $n \longrightarrow \infty$, by Kolmogorov's strong law of large numbers (our theorem with $r = 1$), this upper bound, divided by n^2 , converges a.s. to

$\int h_m(x) dF(x)$, which is arbitrarily small for m large enough. Hence

$$(30) \quad n^{-2} \sum_{1 \leq i < j \leq n} [f^{(j)}(X_i, *) + f^{(j)}(X_j, *)] \longrightarrow 0 \quad \text{a.s.}$$

Finally, consider the first sum in (27),

$$(31) \quad T_n = \sum_{1 \leq i < j \leq n} f_2^{(j)}(X_i, X_j) .$$

By (26), $f_2^{(j)}(x, y) = f_2^{(j)}(y, x)$ and

$$(32) \quad \int f_2^{(j)'}(x, y) dF(y) = 0 .$$

Hence

$$(33) \quad \begin{aligned} E [T_n^2] &= \sum_{1 \leq i < j \leq n} \iint [f_2^{(j)}(x, y)]^2 dF(x) dF(y) \\ &= \sum_{j=2}^n (j-1) \iint [f_2^{(j)}(x, y)]^2 dF(x) dF(y) . \end{aligned}$$

It follows from (26) that

$$\begin{aligned} \iint [f_2^{(j)}(x, y)]^2 dF(x) dF(y) &\leq \iint [f^{(j)}(x, y)]^2 dF(x) dF(y) \\ &= \sum_{v=1}^{j^2} \iint_{v-1 < |f(x, y)| \leq v} [f(x, y)]^2 dF(x) dF(y) \leq \sum_{v=1}^{j^2} v a_v . \end{aligned}$$

Hence

$$(34) \quad E [T_n^2] \leq \sum_{j=2}^n (j-1) \sum_{v=1}^{j^2} v a_v \leq n^2 \sum_{v=1}^{n^2} v a_v .$$

Now, by (31) and (32), $E [T_n | X_1, \dots, X_m] = T_m$ for $m < n$. We apply Lemma 2 with $Z_n = T_n$, $a = b = 2$. We have

$$\sum_{k=1}^{\infty} 2^{-4k} E [T_{2^k}^2] \leq \sum_{k=1}^{\infty} 2^{-2k} \sum_{v=1}^{2^{2k}} v a_v = \sum_{v=1}^{\infty} \sum_{\substack{2^{2k} > v \\ k \geq 1}} 2^{-2k} v a_v < 2 \sum_{v=1}^{\infty} a_v < \infty .$$

Hence

$$(35) \quad n^{-2} T_n \longrightarrow 0 \quad \text{a.s.}$$

It now follows from (27), (28), (30), (31) and (35) that $n^{-2} S'_n \longrightarrow 0$ a.s. By Lemma 3 this implies $n^{-2} S_n \longrightarrow 0$ a.s. The proof of the theorem for $r = 2$ is complete.

REFERENCES

- [1] J. L. Doob, Stochastic Processes, John Wiley and Sons, New York, 1953.
- [2] Wassily Hoeffding, "A class of statistics with asymptotically normal distribution," Ann. Math. Stat., Vol. 19 (1948), pp. 293-325.
- [3] Pranab Kumar Sen, "On some convergence properties of U-statistics," Calcutta Stat. Assn. Bull., Vol. 10 (1960), pp. 1-18.