

TESTS CONCERNING RANDOM POINTS ON A CIRCLE

BY

NICOLAAS H. KUIPER

(Communicated by Prof. H. FREUDENTHAL at the meeting of October 31, 1959)

1. *Introduction*

H. KLOMP, Professor of zoology at Wageningen, suggested the following problem to me:

Let $M = (\varphi_1, \dots, \varphi_m)$ be a sequence of angles ($0 \leq \varphi_i < 2\pi$) representing compass-directions into which m birds, or groups of birds, have been seen flying on migration at a given place on earth at a given time, for example at a given day. M is assumed to be a sequence of m independent observations concerning an unknown random direction φ . Let $M' = (\varphi'_1, \dots, \varphi'_n)$ be analogous with respect to φ' , concerning a second place-time.

1. How to test the nullhypothesis $\varphi \cong \varphi'$?¹⁾
2. How to define and estimate the expectation of a random direction? How to test whether these expectations coincide for φ and φ' ? Observe that in general the value $E(\varphi)$ is not significant in this respect.
3. How to define and estimate a degree of the birds preference for a particular compass-direction, the so-called degree of the orientation of the birds?
4. Given M , how to test the nullhypothesis that the birds have no preference in direction at all, or that they fly according to a given (theoretical) random direction φ .

In this paper we deal with these problems. Different circumstances may invite to different approaches. We first consider in this § some obvious methods. In the other §§ we study a non-parametric test.

In a *first approach* we assume that a strong inclination of the birds is present, so that after a suitable choice of the "coordinate" φ , that is of the direction to be called $\varphi = 0$, it can be assumed that the greater part of the probability-mass is concentrated in a relatively small φ -interval far away from $\varphi = 0$ and from $\varphi = 2\pi$. This φ -interval is then *considered as part of* $-\infty < \varphi < \infty$. Standard methods can be applied now. The exp-

¹⁾ The symbol \cong means "Having the same cumulative distribution function". Compare [5].

tation of the direction φ will be defined and given by $E(\varphi)$ and an estimate is $\frac{1}{m} \sum_1^m \varphi_i$. A measure for the degree of the orientation is

$$\{E(\varphi - E\varphi)^2\}^{-1}$$

which can also be estimated in the usual way. For m sufficiently large, confidence intervals based on normal approximations can be given, and results on different place-times can be compared.

If one has to compare several place-times it is clarifying to represent the corresponding confidence intervals for $E(\varphi)$ and $[E(\varphi - E\varphi)^2]^{\frac{1}{2}}$ in a plane with these numbers as polar coordinates.

WESTENBERG [6] applied related methods also under weaker conditions, in which he indicated a choice for the direction which shall have coordinate $\varphi=0$. Personally we hesitate to go as far as he does.

A *second approach* consists in dividing the interval $0 \leq \varphi < 2\pi$ in sub-intervals, and applying the goodness-of-fit-*Chi-square method*. This may be a good method in practice, for example in the mentioned birds-case it was, because the directions were grouped from the beginning in 32 intervals.

In a *third approach* one considers the random euclidean unit-vector $\mathbf{z} = (\mathbf{x}, \mathbf{y}) = (\cos \varphi, \sin \varphi)$ as a *two-dimensional random vector*. The "expectation of the direction" can be defined as the unit vector $E(\mathbf{z}) / \sqrt{E(\mathbf{z})^2}$ in case $E(\mathbf{z}) \neq 0$. If $E(\mathbf{z}) = 0$ the preferences for different directions cancel. If $E(\mathbf{z})$ has length one, then $P(\mathbf{z} = E(\mathbf{z})) = 1$ so that the degree of orientation is maximal.

A measure for the degree of orientation is $(E(\mathbf{z}))^2$, which obeys $0 \leq (E(\mathbf{z}))^2 \leq 1$. The vector $E(\mathbf{z})$ then comprises both interesting parameters as its direction and length.

If χ_2 is the circle-symmetric standard two-dimensional normal random vector with density $(2\pi)^{-1} \exp. -\frac{1}{2}(x^2 + y^2)$ and σ is such an automorphism of the two-dimensional vector space that any linear function with argument \mathbf{z} has the first and second moments in common with the same linear function with argument $E(\mathbf{z}) + \sigma\chi_2$, then if $\mathbf{z}_1, \dots, \mathbf{z}_m$ are m independent replicates of \mathbf{z} and if m is *sufficiently large*, the random vector $\bar{\mathbf{z}} = \frac{1}{m} \sum_{i=1}^m \mathbf{z}_i$ can, according to the central limit theorem, be approximated by the normal random vector:

$$E(\mathbf{z}) + \frac{1}{\sqrt{m}} \sigma \chi_2.$$

In particular if \mathbf{z} has uniform distribution on the unit circle, one finds that $E(\mathbf{z}) = 0$, σ is a scalar, and

$$\sigma^2 = E(\cos^2 \varphi) = \int_0^{2\pi} \cos^2 \varphi \cdot \frac{d\varphi}{2\pi} = \frac{1}{2}.$$

Hence in this case $\bar{\mathbf{z}}$ is approximated by the random vector

$$\frac{1}{\sqrt{2m}} \boldsymbol{\chi}_2,$$

and we may test the null hypothesis that \mathbf{z} has uniform distribution on the unit circle, with the statistic ¹⁾:

$$\text{approximation: } \left(\frac{\mathbf{z}_1 + \dots + \mathbf{z}_m}{m} \right)^2 = (\bar{\mathbf{z}})^2 \cong \frac{1}{2m} \boldsymbol{\chi}_2^2.$$

For any distribution of \mathbf{z} on the unit circle the second moment of any linear function $f(\mathbf{z})$ with gradient 1, that is a function of the kind $f(\mathbf{z}) = \mathbf{x} \cos \alpha + \mathbf{y} \sin \alpha$, α constant, is immediately seen to be ≤ 1 .

Hence also the variance of any such function is ≤ 1 . This has as a consequence that the linear transformation σ has the property that for any vector z , length $\sigma z \leq$ length z .

Again we assume m sufficiently large so that we could approximate $\bar{\mathbf{z}}$ with a normal random vector. If we use, under these assumptions, a circular confidence region with confidence level α for $E(\mathbf{z})$, obtained by putting

$$(\bar{\mathbf{z}} - E(\mathbf{z}))^2 \cong \frac{1}{m} \boldsymbol{\chi}_2^2 \quad (\text{which is not true})$$

then the true confidence level is certainly at most α .

However, this method is very *inefficient* in case $(E(\mathbf{z}))^2$ is close to 1, that is in case the degree of orientation is large.

§ 2. The non-parametric statistic \mathbf{V}_n

The main aim of this paper is a fourth approach to the problem. This approach, which could be called non-parametric, is connected with the *Kolmogorov* and *Kolmogorov-Smyrnov* tests. These tests are modified so that they can be applied to random points on a circle instead of on a line.

Instead of $\boldsymbol{\varphi}$ in § 1, we now use the random variable

$$\mathbf{x} = \frac{1}{2\pi} \boldsymbol{\varphi}, \quad 0 \leq x < 1.$$

In the sequel we use the notations of KUIPER [5]. The residu class modulo 1 of a real number x or a set of numbers W is denoted by \bar{x} or \bar{W} respectively. The set of residu classes mod 1, with *coordinates* x , $0 \leq x < 1$, is the circle to be considered. The cumulative frequency (c.fr.) function of a given set W of points x_1, \dots, x_n , $0 \leq x_j < 1$, and the cumulative distribution (c.d.) function of a random element \mathbf{y} , $0 \leq \mathbf{y} < 1$, starting cumulating at the value b and jumping from $x=1$ to $x=0$, are

$$F_{\bar{W}}^b(x) \text{ and } F^b(x), \text{ whereas} \\ D_{\bar{W}}^b(x) = F_{\bar{W}}^b(x) - F^b(x), \quad D_{\bar{W}}^0(x) = D_W(x).$$

¹⁾ An example of an alternative against which this test does *not* hold, is the case that the probability mass is equally divided over the vertices of a regular polygon.

We assume that $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent random variables isomorous with a random variable \mathbf{x} with values in $[0, 1)$. And we want to give a test for the nullhypothesis

$$(2.1) \quad H_0 : \mathbf{x} \cong \mathbf{y}.$$

We obtained in [5].

$$(2.2) \quad D_W^b(x) - D_W(x) \text{ is constant.}$$

Consider $\sup_x D_W^b(x)$ and $\inf_x D_W^b(x)$. For $b=0$ these values are denoted in the literature by D_n^+ and $-D_n^-$. KOLMOGOROV uses such statistics and also $\max(D_n^+, D_n^-)$ in order to test H_0 in view of one-sided or two-sided alternatives. (In the general case of real random variables, cumulating starts at $-\infty$).

Analogously one might suggest $\sup_x D_W^b(x)$ and/or $\inf_x D_W^b(x)$ for tests concerning distributions on the circle of reals mod 1. However, the values of these statistics depend on the value b of x at which we start cumulating!

From (2.2) follows that

$$V_W^b = \sup_x D_W^b(x) - \inf_x D_W^b(x)$$

is a function of the sequence $W = (x_1, \dots, x_n)$, $0 \leq x_i < 1$, which is independent of b , and so we may substitute just as well $b=0$ (and omit b in the notation).

$$(2.3) \quad V_W = \sup_x D_W(x) - \inf_x D_W(x).$$

Instead of V_W we occasionally write V_n .

From a random set \mathbf{W} we obtain the random variable:

$$V_{\mathbf{W}} = \mathbf{V}_n = \sup_x D_{\mathbf{W}}(x) - \inf_x D_{\mathbf{W}}(x) = \mathbf{D}_n^+ + \mathbf{D}_n^-$$

which can be used on the circle by passing to the reals mod 1, and which, assuming the nullhypothesis, is independent of the c.d. function on this circle and independent of the point at which we start cumulating.

If b is a constant then we will use the symbol b also for the random variable which has a probability one of taking the value b . If \mathbf{z}_n is a sequence of random variables and the c.d.-functions of \mathbf{z}_n converge for each value to that of a random variable \mathbf{z} , then we say that the limit* of \mathbf{z}_n for $n \rightarrow \infty$ is isomorous with \mathbf{z} :¹⁾

$$\lim_{n \rightarrow \infty}^* \mathbf{z}_n \cong \mathbf{z}.$$

Now, as it is well known that

$$\lim_{n \rightarrow \infty}^* \mathbf{D}_n^+ \cong \lim_{n \rightarrow \infty}^* \mathbf{D}_n^- \cong 0$$

also

$$\lim_{n \rightarrow \infty}^* \mathbf{V}_n \cong 0.$$

¹⁾ This limit* should not be confused with the customary limit of a random infinite sequence.

Consequently we suggest the statistic $V_W = V_n$ in (2.3) for testing whether n given points are independent values of a given theoretical random point on the circle.

3. An asymptotic formula for the c.d. function of \mathbf{V}_n

An asymptotic formula for the distribution-function of \mathbf{V}_n can be obtained from a result of D. DARLING [2]. DARLING has, assuming the null hypothesis

$$(3.1) \quad P(\text{for all } x: -a < \sqrt{n} \mathbf{D}_n(x) < b) = \Phi_n(a, b), \quad a \geq 0, \quad b \geq 0,$$

where

$$\Phi_n(a, b) = \Phi(a, b) + \frac{1}{6\sqrt{n}} \left(\frac{\partial}{\partial a} + \frac{\partial}{\partial b} \right) \Phi(a, b) + o\left(\frac{1}{\sqrt{n}}\right),$$

and

$$\Phi(a, b) = \sum_{j=-\infty}^{\infty} \{e^{-2j^2(a+b)^2} - e^{-2(ja+(j-1)b)^2}\}.$$

The following computation is due to DARLING and the author. The density of the random vector

$$(\sqrt{n} \mathbf{D}_n^-, \sqrt{n} \mathbf{D}_n^+) = (\mathbf{a}, \mathbf{b}) = (-\inf \sqrt{n} \mathbf{D}_W(x), \sup \sqrt{n} \mathbf{D}_W(x))$$

is

$$\frac{\partial^2 \Phi_n}{\partial a \partial b}$$

Hence

$$(3.2) \quad \left\{ \begin{aligned} P(\sqrt{n} \mathbf{V}_n \leq c) &= P(\mathbf{b} - (-\mathbf{a}) < c) = P(\mathbf{a} + \mathbf{b} < c) \\ &= \int_{b=0}^c \left\{ \int_{a=0}^{c-b} \frac{\partial^2 \Phi_n}{\partial a \partial b} da \right\} db \\ &= \int_{b=0}^c \left\{ \frac{\partial}{\partial b} \Phi_n(a, b) \right\}_{a=c-b} db \\ &= A(c) + \frac{B(c)}{\sqrt{n}} + o\left(\frac{1}{\sqrt{n}}\right). \end{aligned} \right.$$

As

$$\left\{ \frac{\partial}{\partial b} \Phi(a, b) \right\}_{a=c-b} = \sum_{j=-\infty}^{\infty} \{-4j^2 c e^{-2j^2 c^2} + 4(j-1)(jc-b) e^{-2(jc-b)^2}\}$$

we have

$$\begin{aligned} A(c) &= \int_{b=0}^c \left\{ \frac{\partial}{\partial b} \Phi(a, b) \right\}_{a=c-b} db = \sum_{j=-\infty}^{\infty} \left\{ -4j^2 c^2 e^{-2j^2 c^2} + (j-1) \int_{b=0}^c e^{-2(jc-b)^2} db \right\} \\ &= \sum_{j=-\infty}^{\infty} \left\{ -4j^2 c^2 e^{-2j^2 c^2} + \right. \\ &\quad \left. + (j-1) (e^{-2(j-1)^2 c^2} - (j-1) e^{-2j^2 c^2}) \right\} \\ &= \sum_{j=-\infty}^{\infty} e^{-2j^2 c^2} \{-4j^2 c^2 + j - (j-1)\} \\ &= \sum_{j=-\infty}^{\infty} (1 - 4j^2 c^2) e^{-2j^2 c^2} \\ &= 1 - \sum_{j=1}^{\infty} 2(4j^2 c^2 - 1) e^{-2j^2 c^2}. \end{aligned}$$

From (3.1) (3.2) one also obtains

$$B(c) = \frac{2}{6} \cdot \frac{d}{dc} A(c) = \frac{8}{3} c \sum_{j=1}^{\infty} j^2 (4j^2 c^2 - 3) e^{-2j^2 c^2}.$$

Hence

$$(3.3) \quad P\{\sqrt{n} \mathbf{V}_n \leq c\} = 1 - \sum_{j=1}^{\infty} 2(4j^2 c^2 - 1) e^{-2j^2 c^2} + \frac{8}{3\sqrt{n}} c \sum_{j=1}^{\infty} j^2 (4j^2 c^2 - 3) e^{-2j^2 c^2} + O\left(\frac{1}{n}\right).$$

For $c > \frac{1}{2}$ a reasonable approximation is obtained from first terms of the series as follows:

$$(3.4) \quad 1 - 2(4c^2 - 1) e^{-2c^2} + \frac{8c}{3\sqrt{n}} (4c^2 - 3) e^{-2c^2}.$$

The formula (3.3) was compared with the result concerning 200 independent samples of 10 numbers between 0 and 1, in three decimal places, obtained from a table of random numbers. The c.fr. function of the 200 values of V_{10} so obtained is given in figure 1. It is seen to be in good agreement with the values according to the formula given in table 1. (The cumulation in the figure goes from right to left.)

TABLE 1

c	$u = c/\sqrt{10}$	$v = P(\mathbf{V}_n > c/\sqrt{10})$
1.0	0.316	0.693
1.1	.348	.528
1.2	.379	.377
1.3	.411	.252
1.4	.443	.158
1.5	.474	.093
1.6	.506	.052
1.7	.536	.027
1.8	.569	.0135
1.9	.600	.0063

§ 4. The statistic $\mathbf{V}_{n,m}$

Analogously one may consider two random c.fr. functions $\mathbf{F}_n^b(x)$ and $\mathbf{F}_m^b(x)$ concerning independent samples of size n and m of two unknown random variables on the circle of reals mod 1 \mathbf{x} and \mathbf{y} , starting the cumulation at b . We want to test the nullhypothesis $\mathbf{x} \cong \mathbf{y}$ in view of "values" $F_n^b(x)$ and $F_m^b(x)$ that the random c.fr. functions $\mathbf{F}_n^b(x)$ and $\mathbf{F}_m^b(x)$ have taken. Let

$$D_{n,m}^b(x) = F_n^b(x) - F_m^b(x).$$

Then $\sup_x D_{n,m}^b(x) - \inf_x D_{n,m}^b(x)$

is independent of the point b , and so we may substitute just as well $b=0$ (and omit b in the notation).

Let

$$V_{n,m} = \sup D_{n,m}(x) - \inf D_{n,m}(x) = D_{n,m}^+ + D_{n,m}^-.$$

The statistic $V_{n,m}$ is independent of the point at which we start cumulating; it is also independent of the c.d. function of $\mathbf{x} (\cong \mathbf{y})$.

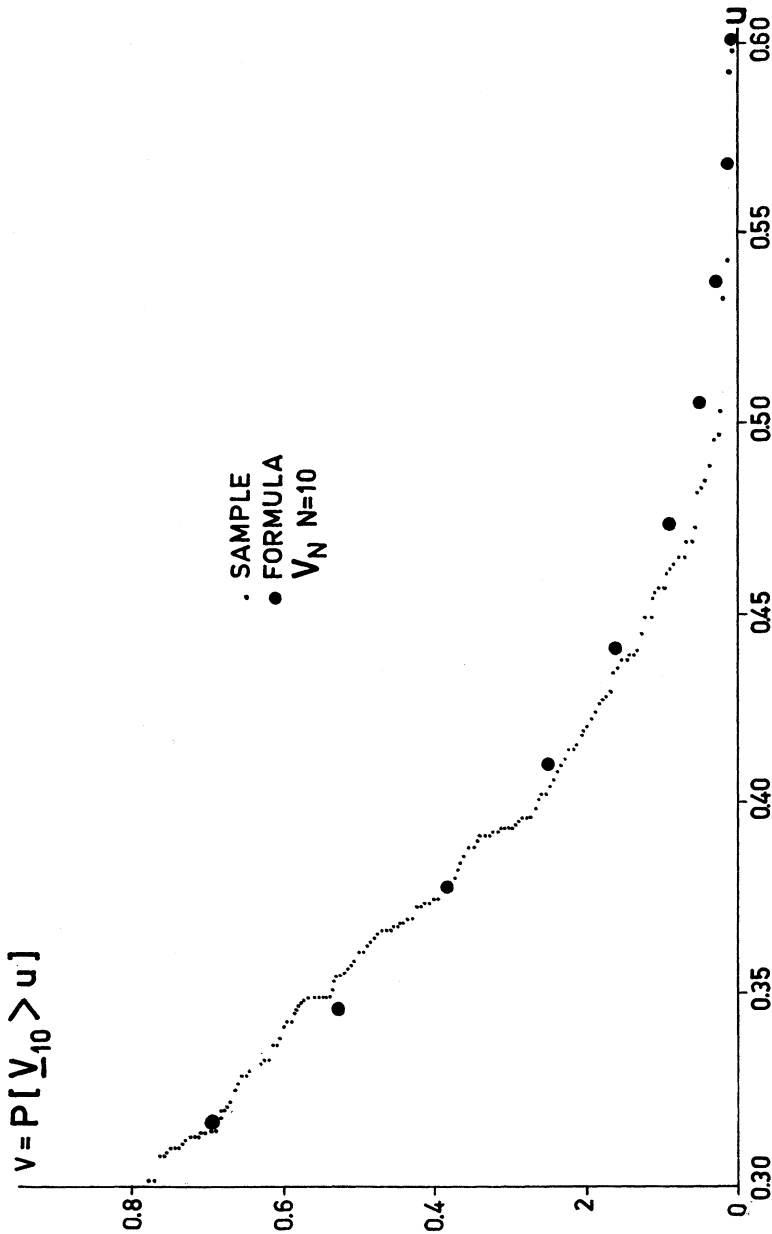


Fig. 1

We remark

$$\lim_{m \rightarrow \infty}^* \mathbf{V}_{n,m} \cong \mathbf{V}_n$$

and recall

$$\lim_{n \rightarrow \infty}^* \mathbf{V}_n \cong 0.$$

We suggest the statistic $\mathbf{V}_{n,m}$ for testing whether in case n given independent values of an unknown random point \mathbf{x} and m given independent values of an unknown random point \mathbf{y} can come from the same continuous distribution $\mathbf{x} \cong \mathbf{y}$ on the circle.

In order to apply the test one needs the c.d. function of $\mathbf{V}_{n,m}$ which, however, we did not yet determine in general. For the case $m=n$ we can obtain this c.d. function as follows from a formula of GNEDENKO [3], recently improved by KEMPERMAN [5]: § 4, formula (9).

$$(4.1) \quad \left\{ \begin{aligned} P_n(a, b) &= P \left(\text{for all } x: -\frac{a}{n} < \mathbf{D}_{n,n}(x) < \frac{b}{n} \right) = \\ &= g_0 + (3g_0 - g_2)/(24n) + \left(\frac{9}{2}g_0 - 3g_2 - \frac{16}{5}g_3 + \frac{1}{2}g_4 \right) / (24n)^2 + O(n^{-3}) \end{aligned} \right.$$

with

$$\begin{aligned} (-1)^r g_r(a, b, n) &= (-1)^r g_r = \sum_{k=-\infty}^{\infty} \left\{ H_{2r}^* \left(\frac{2kc}{\sqrt{2n}} \right) - H_{2r}^* \left(\frac{2a+2kc}{\sqrt{2n}} \right) \right\} \\ H_{2r}^*(x) &= \left(\frac{d}{dx} \right)^{2r} e^{-x^2/2}, \quad c = a + b. \end{aligned}$$

We will use different formulas obtained as follows:

Let

$$(4.2) \quad \Psi_n(a, b) = P_n(a\sqrt{n}, b\sqrt{n}) = P \left(-\frac{a}{\sqrt{n}} < \mathbf{D}_{n,n}(x) < \frac{b}{\sqrt{n}} \text{ for all } x \right)$$

and

$$g_r(a\sqrt{n}, b\sqrt{n}, n) = h_r(a, b, n) = h_r.$$

Then

$$(-1)^r h_r = \sum_{k=-\infty}^{\infty} \left\{ H_{2r}^*(\sqrt{2}kc) - H_{2r}^*(\sqrt{2}(a+kc)) \right\}$$

and

$$(4.3) \quad \left\{ \begin{aligned} \Psi_n(a, b) &= h_0 + (3h_0 - h_2)/(24n) + \\ &+ \left(\frac{9}{2}h_0 - 3h_2 - \frac{16}{5}h_3 + \frac{1}{2}h_4 \right) / (24n)^2 + O(n^{-3}). \end{aligned} \right.$$

For our purpose it will be sufficient to use

$$(4.4) \quad \Psi_n(a, b) = h_0 + (3h_0 - h_2)/(24n) + O(n^{-2})$$

with

$$h_0 = \sum_{j=-\infty}^{\infty} (e^{-j^2 c^2} - e^{-(a+jc)^2}), \quad c = a+b$$

$$h_2 = \sum_{j=-\infty}^{\infty} (2j^2 c^2 - 1) e^{j^2 c^2} - (2(a+jc)^2 - 1) e^{-(a+jc)^2}.$$

The density of the random vector

$$(\mathbf{a}, \mathbf{b}) = (-\inf_x \sqrt{n} \mathbf{D}_{n,n}(x), \sup_x \sqrt{n} \mathbf{D}_{n,n}(x))$$

is therefore

$$\frac{\partial^2 \Psi_n(\mathbf{a}, \mathbf{b})}{\partial \mathbf{a} \partial \mathbf{b}}.$$

Hence

$$(4.2) \quad \left\{ \begin{aligned} P(\sqrt{n} \mathbf{V}_{n,n} \leq c) &= P(\mathbf{a} + \mathbf{b} \leq c) = \int_{b=0}^c \left\{ \int_{a=0}^{c-b} \frac{\partial^2 \Psi_n(\mathbf{a}, \mathbf{b})}{\partial \mathbf{a} \partial \mathbf{b}} d\mathbf{a} \right\} d\mathbf{b} \\ &= \int_{b=0}^c \left\{ \frac{\partial}{\partial \mathbf{b}} \Psi_n(\mathbf{a}, \mathbf{b}) \right\}_{\mathbf{a}=c-b} d\mathbf{b}. \end{aligned} \right.$$

The computations are as in § 3. One obtains:

$$(4.3) \quad \left\{ \begin{aligned} P(\sqrt{n} \mathbf{V}_{n,n} \leq c) &= 1 - \sum_{j=1}^{\infty} 2(2j^2 c^2 - 1) e^{-j^2 c^2} + \\ &+ \frac{1}{6n} \left(1 + \sum_{j=1}^{\infty} j^2 c^2 (2j^2 c^2 - 7) e^{-j^2 c^2} \right) + 0(n^{-2}). \end{aligned} \right.$$

Some critical regions concerning \mathbf{V}_n and $\mathbf{V}_{n,n}$ are given in tables 2 and 3. Conclusions concerning $\mathbf{V}_{n,m}$ for $n < m$ can occasionally be obtained from these tables and the fact that if $n < m$,

$$P(\sqrt{n} \mathbf{V}_n > c) < P(\sqrt{n} \mathbf{V}_{n,m} > c) < P(\sqrt{m} \mathbf{V}_{n,m} > c).$$

TABLE 2
Critical regions for the \mathbf{V}_n -test. $P(\sqrt{n} \mathbf{V}_n > c) = \alpha$

$\alpha \backslash n$	10	20	30	40	100	∞
.10	1.4877	1.5322	1.5503	1.5608	1.5839	1.6196
.05	1.6066	1.6564	1.6760	1.6869	1.7110	1.7473
.01	1.8391	1.9027	1.9153	1.9375	1.9637	2.0010

TABLE 3
Critical regions for the $\mathbf{V}_{n,n}$ -test. $(P(\sqrt{n} \mathbf{V}_{n,n} > c) = \alpha$

$\alpha \backslash n$	10	20	30	40	100	∞
.10	2.2429	2.2663	2.2743	2.2783	2.2855	2.2905
.05	2.4041	2.4376	2.4488	2.4543	2.4643	2.4710
.01	2.6125	2.6988	2.7352	2.7556	2.7974	2.8298

Landbouwhogeschool Wageningen

REFERENCES

1. BIRNBAUM, Z. W. and R. PYKE, On some distributions related to the statistic D_n^+ . (Ann. of Math. Stat. 29, 179-187 (1958).
2. DARLING, D. A., To appear in the Bulletin of the Amer. Math. Soc.
3. GNEDENKO, B. V., Some results on the maximal deviation between two empirical distributions. Dokl. Akad. Nauk SSSR, 82, 661-663 (1952).
4. KEMPERMAN, J. H. B., Some exact formulae for the Kolmogorov-Smyrnov distributions. Indag. Math. XIX, 535-540. = Proc. Amsterdam (1957).
5. ———, Asymptotic expansions for the Smirnov test and the range of cumulative sums. Ann. of Math. Stat. 30, 448-463 (1959).
6. KUIPER, N. H., Annals of Math. Stat. 30, 251-252 (1959).
7. ———, On the random cumulative frequency function. This Proc. 32 Amsterdam (1959)
8. WESTENBERG, J., The median and interquartile range test applied to frequency distributions plotted on a circular axis. Indag. Math. XII, 378-381 = Proc. Amsterdam (1950).