

# ON SERIAL NUMBERS

By E. J. GUMBEL

*New School for Social Research*

In this paper we consider a continuous variate and unclassified observations. It is well known that there are two step functions, which we may trace for a given series of observations. We will show that the differences between the two ways of plotting play an important rôle for certain graphical methods used by engineers.

To obtain one and only one series of observations we adjust the cumulative frequencies. The corrections thus introduced depend upon the theoretical distribution which is adequate for the observations. Later we deal with the relation between serial numbers and grades. Finally we construct confidence bands for the comparison between theory and observations.

**1. Theory and observations.** If we arrange  $n$  observations in order of increasing magnitude, and write each as often as it occurs, there will be a first,  $x_1$ , the smallest value, a second,  $x_2$ , an  $m$ th,  $x_m$  the penultimate,  $x_{n-1}$ , and the last,  $x_n$ , i.e., the greatest value. The index  $m$  is called the observed cumulative frequency, or simply the rank. It is usual to draw the observations  $x_m$  along the abscissa, and the rank  $m$  along the ordinate. The step function starts with a vertical line from the value  $x_1$  of the abscissa to the point with the coordinates  $1, x_1$ , and in general consists of the horizontal lines from the point  $m, x_m$  to the point  $m, x_{m+1}$  and the vertical lines from the point  $m, x_{m+1}$  to the point  $m+1, x_{m+1}$ . The step function ends with the point  $n, x_n$ . We call this graph the step function  $(m, x_m)$ . However, another step function which is derived from the observations arranged in decreasing magnitude is equally legitimate. This step function starts from the point with the coordinates  $0, x_1$ , and in general consists of the horizontal lines from the point  $m-1, x_m$  to  $m-1, x_{m+1}$  and the vertical lines from the point  $m-1, x_{m+1}$  to the point  $m, x_{m+1}$  and ends with the point  $n-1, x_n$ . We call it the step function  $(m-1, x_m)$ . Let  $F(x)$  be the probability of a value equal to or less than  $x$ . Then the continuous theoretical curve, the ogive, which we compare to the step functions is  $nF(x), x$ . The question is whether we have to use the step function  $(m, x_m)$  or the step function  $(m-1, x_m)$ .

The differences between the two ways of plotting are rarely mentioned in the statistical literature. If we plot instead of the rank  $m$  the reduced rank  $m/n$ , the differences between the two ways of plotting are of the order  $1/n$ . It is generally tacitly assumed that this difference may be neglected. This will not hold if  $n$  is small.

In the following we show two other ways of plotting the observations where the differences between the two observed curves play an important role. Sup-

pose that the probability  $F(x)$  and the density of probability,  $f(x)$ , henceforth called the distribution are such that it is possible to introduce a reduced variate

$$(1) \quad z = \frac{x - a}{b},$$

which has no dimension. In general, the constant  $a$  will be a certain mean, and the constant  $b$  a certain measure of dispersion. Furthermore, the constants may be linear functions of these characteristics. Neither the probability  $G(z)$  of a value equal to or less than  $z$

$$(2) \quad G(z) = F(x),$$

nor the reduced distribution

$$(3) \quad g(z) = bf(a + bz)$$

contain constants. The *equiprobability test* consists in the following procedure: We attribute to the  $m$ th observation  $x_m$  the relative frequency  $m/n$ , and determine from a probability table a value  $z$ , such that

$$(4) \quad G(z) = m/n.$$

The variate  $x$  is plotted on the ordinate, and the reduced variate  $z$  on the abscissa. Then the points  $x_m, z$  must be situated close to the straight line (1). To apply this comparison between theory and observations, we need not even calculate the constants. For the normal distribution the application of this test is greatly facilitated by the use of probability paper.

The difficulty is that we may as well choose the frequency

$$(4') \quad G(z) = (m - 1)/n,$$

and determine the corresponding values of  $z$ . Therefore, we have two lines (1) instead of one. The difference between the two series will be large for the first and last few observations. For the first series the last observation cannot be plotted on probability paper; for the second series the first observation cannot be plotted.

The same difficulty exists for the "return period." If the observations of a continuous variate are made at regular intervals in time which are taken as units, we may as in [4] define the theoretical return periods  $T(x)$  of a value equal to or greater than  $x$  as

$$(5) \quad T(x) = \frac{1}{1 - F(x)}.$$

The comparison of the theoretical with the observed return periods gives a test for the validity of a theory. However, there are two series of observations, namely, the exceedance intervals

$$(6) \quad 'T(x_m) = \frac{n}{n-m}; \quad m = 1, 2 \dots n-1$$

and the recurrence intervals

$$(7) \quad ''T(x_m) = \frac{n}{n-m+1}; \quad m = 1, 2 \dots n.$$

The two expressions (6) and (7) differ widely for the high ranks. The penultimate observation, for example, has an exceedance interval  $n$ , whereas the recurrence interval is only  $n/2$ . This contradiction and the difficulty arising for the equiprobability test show that the question of choosing the observed cumulative frequency of the  $m$ th observation has a practical significance.

The equiprobability test and the comparison between the observed and the theoretical return period may be combined on probability paper. The variate  $x$  is plotted on the ordinate, the reduced variate  $y$  on the abscissa. But instead of  $y$  we write the probability  $F(x)$  and the return period  $T(x)$ . If the theory holds, the observations must be scattered around the straight line (1).

But all these methods presuppose that we know whether we have to attribute to  $x_m$  the rank  $m$  or the rank  $m-1$ . Sometimes a compromise has been proposed which consists in attributing to  $x_m$  neither  $m$  nor  $m-1$ , but the arithmetic mean of both,  $m - \frac{1}{2}$ . In other words, the index  $m$  is no longer considered to be an integer. In such cases, we call  $m$  the *serial number*.

The corrected frequency  $m - \frac{1}{2}$  may be accepted for the comparison between the step function and the probability curve. However, for the return period and for the equiprobability test this method leads to serious difficulties. The corrected return periods, which have been proposed by Hazen [7] and have been used by M. Kimball [8] are

$$(6) \quad T(x_m) = \frac{n}{n-m+1/2}.$$

The last among  $n$  observations has a return period  $2n$ . This idea does not seem to be sound. No statistical device can increase the number of observations beyond  $n$ .

**2. The adjusted frequency of the  $m$ th observation.** The use of  $m$ ,  $m-1$ , or  $m - \frac{1}{2}$  as frequency of the  $m$ th observations amounts to considering the  $m$ th value as being fixed. To obtain one and only one step function we consider  $x_m$  as a statistical variate. This will lead to the determination of the most probable serial number and of the corresponding probability as a function of  $m$  and  $n$ .

The  $m$ th observation is such that there are  $m-1$  observations below it and  $n-m$  observations above it. Consequently, the distribution  $w_n(x, m)$  of the  $m$ th observation is

$$(9) \quad w_n(x, m) = \binom{n}{m} m [F(x)]^{m-1} [1 - F(x)]^{n-m} f(x).$$

The variate  $x_m$  is simply called  $x$  as each value of  $x$  has a certain density of probability of being the  $m$ th. To distinguish between  $(x)$  and  $w_n(x, m)$ , the first distribution is referred to as the *initial* distribution. For some simple initial distributions it is possible to calculate exactly the mean and the standard error of the distribution (9). This has been done by Karl Pearson [10] for the normal, the uniform, the exponential, and other skew distributions. The results are very complicated, and do not allow any immediate practical applications.

In the following we determine therefore instead of the mean the mode of the  $m$ th value. The most probable  $m$ th value for which we write  $\tilde{x}_m$  is the solution of

$$\frac{d \log w_n(x, m)}{dx} = 0.$$

We obtain from (9)

$$(10) \quad \frac{m-1}{F(\tilde{x}_m)} f(\tilde{x}_m) - \frac{n-m}{1-F(\tilde{x}_m)} f(\tilde{x}_m) = -\frac{f'(\tilde{x}_m)}{f(\tilde{x}_m)}.$$

In this equation  $m$  is counted in order of increasing magnitude. If we choose the inverse order we obtain the same equation, if we replace the index  $m$  by  $n-m+1$ . Therefore the following results are independent of the order of counting  $m$ .

Equation (10) gives the most probable value  $\tilde{x}_m$  as a function of  $m$  and  $n$ . The function depends upon the distribution.

A rough, first trial solution of (10) may be obtained if we confine our interest to values where neither  $m$  nor  $n-m$  is small in comparison to  $n$ , that is, values which are not extreme. Suppose  $m$  to be of the order  $n/2$ . For increasing numbers of observations, the expression on the left side of (10) become large compared to the right side provided the derivative remains finite, as is generally the case. If we neglect the right-hand member,  $\tilde{x}_m$  is the solution of

$$(11) \quad F(\tilde{x}_m) = \frac{m-1}{n-1}.$$

This expression holds for the uniform distribution where  $f'(x) = 0$ .

The following exact solution is valid for any number of observations and any serial number. Equation (10) will be used in two different ways: First, we suppose  $m$  to be known, we determine the probability  $F(\tilde{x}_m)$  of the most probable  $m$ th value as a function of  $m$  and  $n$ , and attribute this probability to the  $m$ th observation  $x_m$ . By doing so, the probability of the most probable  $m$ th value becomes the *adjusted frequency* of the  $m$ th observation. This leads to one and only one series of observations, and settles our initial question. Later, in section 3, we suppose  $F(\tilde{x}_m)$  to be known, and compute the corresponding most probable  $m$ th observation. This leads to an estimate of the grades (or partition values) from the serial numbers.

To obtain  $F(\tilde{x}_m)$  from (10) we introduce an expression  $\sigma^2(x_m)$  by stating

$$(12) \quad [\sigma^2(x_m)n] = F(\tilde{x}_m)[1 - F(x_m)]f^{-2}(\tilde{x}_m).$$

The brackets are meant to indicate that the product on the left side does not depend upon  $n$ . We shall show later that  $\sigma^2(x_m)$  is under certain conditions the variance of the  $m$ th observation. For the present purpose however this significance is not required. Multiplication of (10) by (12) leads to

$$(13) \quad m - 1 + F(\tilde{x}_m) - nF(\tilde{x}_m) = -f'(\tilde{x}_m)[\sigma^2(x_m)n],$$

or

$$(14) \quad F(\tilde{x}_m) = \frac{m - 1}{n - 1} + \frac{f'(\tilde{x}_m)[\sigma^2(x_m)n]}{n - 1}.$$

The adjusted frequency in (14) is similar to (11). Another expression for the adjusted frequency, derived from (13), is

$$(15) \quad F(\tilde{x}_m) = \frac{m - \frac{1}{2}}{n} + \frac{1}{n} (F(\tilde{x}_m) - \frac{1}{2} + f'(\tilde{x}_m)[\sigma^2(x_m)n]).$$

The adjusted frequency is the compromise  $\frac{m - \frac{1}{2}}{n}$  plus an expression

$$(16) \quad \frac{D}{n} = \frac{1}{n} (F(\tilde{x}_m) - \frac{1}{2} + f'(\tilde{x}_m)[\sigma^2(x_m)n]).$$

The correction,  $D$ , defined by (16) depends upon the initial distribution and has no dimension. In general, it will depend upon the constants which exist in the distribution. If the distribution  $f(x)$  may be written in a reduced form (3), the correction<sup>1</sup>

$$(17) \quad D = G(z) - \frac{1}{2} + g'(z)[\sigma^2(z)n]$$

depends only upon the dimensionless reduced variate  $z$ . For a given initial distribution we choose numerical values for the probability  $G(z) = F(\tilde{x}_m)$  calculate  $g'(z)$  and

$$(18) \quad [\sigma^2(z)n] = \frac{G(z)(1 - G(z))}{g^2(z)}.$$

From (16) we compute a table for the corrections  $D$  as a function of the adjusted frequencies  $F(\tilde{x}_m)$  and obtain for given  $n$  the serial number  $m$  as a function of the adjusted frequencies by

$$(19) \quad m = nF(\tilde{x}_m) + \frac{1}{2} - D.$$

These serial numbers will not be integers. The adjusted frequency  $F(\tilde{x}_m)$  for the  $m$ th observation will then be obtained by linear interpolation.

<sup>1</sup> In previous articles [3, 6] we started from another interpretation of the corrected frequencies and obtained slightly different corrections.

The value and the sign of the correction  $D$  depends upon the distribution. For the asymmetrical exponential distribution, for example, the correction

$$(19') \quad D = -\frac{1}{2},$$

is independent of the variate. This means that we have to use exclusively the step function  $(m-1, x_m)$  as being the best way of plotting. The observed adjusted return periods are the recurrence intervals.

For a symmetrical reduced distribution we have

$$(20) \quad 1 - G(-z) = G(z); \quad g(-z) = g(z); \quad g'(-z) = -g'(z).$$

Therefore, the reduced correction will be

$$(21) \quad D(-z) = -D(z).$$

For the two reduced values  $z$  and  $-z$  of a symmetrical variate the corrections have the same size, but different signs.

A relation similar to (21) holds for two asymmetrical reduced distributions  $g_1(z)$  and  ${}_1g(z)$ , which are symmetrical one to another in the sense

$$(22) \quad G_1(z) = 1 - {}_1G(-z); \quad g_1(z) = {}_1g(-z); \quad g'_1(z) = -{}_1g'(-z).$$

Then, the corrections are

$$(23) \quad D_1(-z) = -{}_1D(z).$$

For any initial distribution  $f(x)$  we read from (19) the adjusted frequency

$$(24) \quad F(\tilde{x}_m) = \frac{m - \frac{1}{2} + D}{n},$$

even for a small number of observations. The question whether to choose  $m/n$  or  $(m-1)/n$  as observed cumulative frequency is settled by (24). We obtain one observed step function, one series for the equiprobability test, and one series of observed return periods

$$(25) \quad T(\tilde{x}_m) = \frac{n}{n - m + \frac{1}{2} - D},$$

which have to be compared to the theoretical continuous curves.

**3. Estimates for the grades.** In the following we use the fundamental formula (15) to determine interesting grades through the  $m$ th values.

We use the term *grade* for the value of a statistical variate which corresponds to a given cumulative probability  $F(x)$  say,  $F(x) = \frac{1}{4}; \frac{1}{2}; \frac{3}{4}$  for quartiles;  $F(x) = \frac{1}{10}, \dots, \frac{9}{10}$  for deciles, and so on. For a given grade, the probability  $F(x)$  the density of probability  $f(x)$  and its derivative are known, and  $m$  is unknown. The value of  $m$  obtained from (15), henceforth called the most probable serial number  $\tilde{m}$ , is the solution of

$$(26) \quad \tilde{m} = nF(x) + 1 - F(x) - f'(x)F(x)(1 - F(x))f^{-2}(x).$$

The corresponding "observed" value  $x_{\tilde{m}}$  is obtained by interpolation between two observed values  $x_{m-1}$  and  $x_m$ , such that

$$m - 1 < \tilde{m} < m.$$

For the median,  $x_0$ , the most probable serial number  $\tilde{m}_0$  is

$$(27) \quad \tilde{m}_0 = \frac{n+1}{2} - \frac{f'(x_0)}{4f^2(x_0)}.$$

The median  $x_0$  itself enters into (27). It has to be eliminated through the condition  $F(x_0) = \frac{1}{2}$ . For the exponential distribution for example we find

$$(27') \quad \tilde{m}_0 = \frac{n}{2} + 1.$$

The most probable serial number of the median for a symmetrical distribution is

$$(28) \quad \tilde{m}_0 = \frac{1}{2}(n+1).$$

This is the usual estimate of the median for any distribution. The estimate obtained from (27) is smaller (larger) than the usual estimate if the median is smaller (larger) than the mode. The difference between the two estimates is due to the fact, that (27) makes use of information about the theoretical distribution whereas this information (if available) is neglected by the usual method.

For symmetrical distributions the most probable serial numbers  $\tilde{m}_1$  and  $\tilde{m}_2$  for two symmetrical grades defined by  $F_1$  and  $F_2 = 1 - F_1$  are according to (26) related by

$$(29) \quad \begin{aligned} \tilde{m}_1 &= nF_1 + 1 - (F_1 + f'_1F_1(1 - F_1)f_1^{-2}) \\ \tilde{m}_2 &= n(1 - F_1) + (F_1 + f'_1F_1(1 - F_1)f_1^{-2}). \end{aligned}$$

The members in brackets have the same size, but opposite signs. Another expression for  $\tilde{m}_2$  is

$$\tilde{m}_2 = (n+1) - [nF_1 + 1 - F_1 - f'_1F_1(1 - F_1)f_1^{-2}]$$

so that, for symmetrical distributions

$$(30) \quad \tilde{m}_1 + \tilde{m}_2 = n + 1.$$

This is to be expected as the  $m$ th value counted upwards is the  $(n - m + 1)$ st value counted downwards.

For the two quartiles  $q_1$  and  $q_2$  the most probable serial numbers  $\tilde{m}(q_1)$  and  $\tilde{m}(q_2)$ , obtained from (29) are

$$(31) \quad \tilde{m}(q_1) = \frac{n+3}{4} - \frac{3 f'(q_1)}{16 f^2(q_1)}; \quad \tilde{m}(q_2) = \frac{3n+1}{4} - \frac{3 f'(q_2)}{16 f^2(q_2)},$$

where  $q_1$  and  $q_2$  have to be eliminated by the use of

$$F(q_1) = \frac{1}{4}; F(q_2) = \frac{3}{4}.$$

For the uniform, the normal and the exponential distribution we obtain the two quartiles from

$$\begin{aligned} \tilde{m}(q_1) &= \frac{n+3}{4} \quad ; \quad \tilde{m}(q_2) = \frac{3n+1}{4} \\ (31') \quad \tilde{m}(q_1) &= \frac{n}{4} + .352; \quad \tilde{m}(q_2) = \frac{3n}{4} + .648 \\ \tilde{m}(q_1) &= \frac{n}{4} + 1 \quad ; \quad \tilde{m}(q_2) = \frac{3n}{4} + 1 \end{aligned}$$

respectively. The last result may also be found from (19') and (24). These estimates differ from the usual estimates by the reason given above.

We now apply the notion of a grade to certain characteristics which are *otherwise* defined. A certain characteristic, say, the mode  $\tilde{x}$  or the mean  $\bar{x}$  have for a given distribution the probabilities  $F(\tilde{x})$  or  $F(\bar{x})$  respectively. These probabilities may be used to define a grade. We determine the corresponding  $m$ th value from (26), and obtain an estimate of the mode or the mean, interpreted as grades by interpolation between the observed  $m$ th values. For a symmetrical distribution these estimates of the mode and mean are identical with the estimates of the median. For an asymmetrical distribution, the most probable serial number  $\tilde{m}(\tilde{x})$  of the mode becomes according to (26)

$$(32) \quad \tilde{m}(\tilde{x}) = (n-1)F(\tilde{x}) + 1.$$

Usually, the mode  $\tilde{x}$  of a continuous variate is estimated by another procedure. The observations are arranged in certain cells. One of them has the largest relative frequency. It will contain the mode. To find its position within the cell, an interpolation formula is applied which reproduces the content of this cell and of the two adjacent cells. By choosing different lengths for the cells and different origins for the classification, the mode can be shifted to the right or to the left. Formula (32) furnishes a determination of the mode from the observations according to the theory, such that the arrangement of the observations into different cells is not needed. Of course, this method can be applied only if we know the theoretical distribution  $f(x)$ . The problem how to estimate the mode is important for distributions where one of the constants may be interpreted as the mode or as a function of the mode [1, 4].

**4. Standard errors of the estimates.** The numerical work involved in the method (26) of estimating the grades is very small. To obtain the standard errors of these estimates we consider the asymptotic properties of the distribution (9). The following results hold therefore only for large numbers of observation. Besides we assume, that the serial number  $m$  is of the order  $n/2$ , i.e. not extreme. It has been shown [2] that under these conditions the distribution



of the  $m$ th value converges toward a normal distribution with a standard error  $\sigma(x_m)$ , where

$$(33) \quad \sigma(x_m)\sqrt{n} = \frac{1}{f(x)} \sqrt{F(x)(1-F(x))}.$$

Although this standard error does not contain  $m$  explicitly, it has a clear meaning for any value of  $x$  as we know from (26), which observation we have to attribute to the probability  $F(x)$ . The classical proof about the approximate normality of the distribution of the median in large samples is a special case of this convergence and the classical standard error of the median,

$$(34) \quad \sigma(x_0)\sqrt{n} = \frac{1}{2f(x_0)},$$

is a special case of (33). The square root in (33) is maximum for  $F(x) = \frac{1}{2}$ . Therefore,

$$(35) \quad \sigma(x_m)\sqrt{n} \leq \frac{1}{2f(x)}.$$

If the variate  $x$  may be reduced through the linear transformation (1) the standard error  $\sigma(z)$  of the reduced variate, called reduced standard error

$$(36) \quad [\sigma(z)\sqrt{n}] = \frac{1}{g(z)} \sqrt{G(z)(1-G(z))},$$

may be calculated as a function of  $z$  where  $z$  corresponds to  $x_m$ . To call attention to the fact that these numerical values do not depend upon  $n$ , they are written in brackets. The standard error of the estimate for  $x_m$  is, according to (2) and (3)

$$(37) \quad \sigma(x_m) = \frac{b}{\sqrt{n}} [\sigma(z)\sqrt{n}].$$

Since the constant  $b$  is a measure of dispersion, the standard error of the estimate of the  $m$ th value is proportional to the standard deviation of the variate. The factors  $b$  and  $1/\sqrt{n}$  show that the standard error of the  $m$ th value is of the same structure as the standard error of the arithmetic mean.

For symmetrical distributions the standard error (33) of the  $m$ th value is also a symmetrical function. The standard errors of the estimate of the two quartiles, and generally of the estimates of two grades defined by  $F$  and  $1 - F$ , are then identical. If the mode coincides with the median, the corresponding standard error of the  $m$ th value is a minimum. For a symmetrical  $U$ -shaped distribution, however, where the density of probability passes through a minimum at the center of symmetry, the median has the largest standard error among the  $m$ th values. An example for such a distribution has been given by Leavens [9]. As the distribution of the  $m$ th value converges towards a normal distribution, it is legitimate to attribute to the mode of the  $m$ th value the standard error (33).

Therefore, for a large number of observations (33) gives the standard error of our estimate of the grades. The standard errors of the estimates (31) of the quartiles are

$$(38) \quad \sigma(q_1)\sqrt{n} = \frac{\sqrt{3}}{4f(q_1)}; \quad \sigma(q_2)\sqrt{n} = \frac{\sqrt{3}}{4f(q_2)}.$$

The arithmetic mean in its usual definition is not an  $m$ th value. Its standard error  $\sigma(\bar{x})$ , where

$$(39) \quad \sigma(\bar{x})\sqrt{n} = \sigma,$$

will, therefore, fall outside of the range of the standard errors of the  $m$ th values. (See graph 1.) If the distribution  $f(x)$  is such that the standard deviation does not exist, it is legitimate to estimate the arithmetic mean as a grade, and calculate it from the corresponding most probable  $m$ th value by introducing  $F(\bar{x})$ ,  $f(\bar{x})$  and  $f'(\bar{x})$  into (26). The standard error of the arithmetic mean interpreted as a grade is

$$(40) \quad \sigma(\bar{x})\sqrt{n} = \frac{1}{f(\bar{x})} \sqrt{F(\bar{x})(1 - F(\bar{x}))}.$$

If we use this estimate of the arithmetic mean for distributions where  $\sigma$  exists, the usual determination of the mean will be more (less) precise than its estimate as a grade if

$$(41) \quad \sigma f(\bar{x}) \leq \sqrt{F(\bar{x})(1 - F(\bar{x}))}.$$

The standard error of the mode estimated as a grade is

$$(42) \quad \sigma(\tilde{x})\sqrt{n} = \frac{1}{f(\tilde{x})} \sqrt{F(\tilde{x})(1 - F(\tilde{x}))}.$$

As the standard error of any characteristic depends upon the way it is estimated from the observations, the standard errors of the mode or mean interpreted as grades differ from the usual standard errors.

**5. The most precise grade.** Equation (33) may be used to define a new grade which has interesting properties. The standard error (33) of the estimate of the  $m$ th value is a function of  $F$ . We ask whether it possesses a minimum (maximum). The corresponding value of the variate,  $\hat{x}$ , may be called *the most (least) precise  $m$ th value* or the most (least) precise grade. To obtain  $\frac{d\sigma(x_m)}{dF}$  it is sufficient to calculate from (33)

$$\frac{nd \log \sigma^2(x_m)}{dx} = \frac{2n\sigma'(x_m)}{\sigma(x_m)}.$$

Therefore the most (least) precise grade is the solution of

$$(43) \quad \frac{f(x)}{F(x)} - \frac{f(x)}{1 - F(x)} - \frac{2f'(x)}{f(x)} = 0.$$

This expression does not vanish if either  $F(x) = \frac{1}{2}$  or  $f'(x) = 0$ . It vanishes if both conditions hold simultaneously. For a symmetrical distribution passing through a mode (minimum), the mode (minimum), estimated as a grade, is the most (least) precise grade. Equation (43) may be written

$$f'(x)f''(x)F(x)(1 - F(x)) = \frac{1}{2} - F(x).$$

If we introduce this expression into (16), we obtain  $D = 0$  and

$$(44) \quad F(\hat{x}) = \frac{m - \frac{1}{2}}{n}.$$

*The most precise  $m$ th value is such that the adjusted frequency is the arithmetic mean of the frequencies  $m/n$  and  $(m - 1)/n$ .*

The most precise  $m$ th value  $\hat{x}$  cannot be calculated from the observations alone. It may be estimated in the same way as any grade by introducing the values  $F(\hat{x})$ ,  $f(\hat{x})$  and  $f'(\hat{x})$  into equation (26).

To show the difference between the most precise grade and the mode we apply the procedure developed above to a skew distribution. The reduced distribution of the largest value  $g(y)$  and the probability  $G(y)$  are

$$(45) \quad g(y) = e^{-y}G(y); \quad G(y) = e^{-e^{-y}}.$$

The relation (1) between the reduced variate for which we write  $y$  instead of  $z$  and the largest value  $x$  is

$$(46) \quad x = u + \frac{y}{\alpha}.$$

where  $u = \tilde{x}$  is the mode and

$$(47) \quad \frac{1}{\alpha} = \frac{\sqrt{6}}{\pi} \sigma.$$

The most probable serial number  $\tilde{m}(u)$  of the mode, obtained from (32) is

$$(48) \quad \tilde{m}(u) = \frac{n + e - 1}{e}.$$

This equation may be used for an estimate of the constant  $u$ .

The reduced variance  $\sigma^2(y)$  obtained from (36) and (45) is

$$(49) \quad (\sigma^2(y)\sqrt{n}) = e^{2y}(e^{e^{-y}} - 1).$$

A table for the reduced standard error  $\sigma(y)\sqrt{n}$  has been given in a previous publication [6]. The value  $\sigma(y)\sqrt{n}$  is plotted in figure 1 for probabilities  $G(y)$  from 0.01 to 0.95. The standard error has a minimum for a value of  $y$  located to the left of the mode  $\tilde{y} = 0$ . On the same graph are plotted the reduced standard errors for the normal distribution. As the normal reduced variate  $z$  differs from the reduced variate  $y$ , two different scales are used for the variates. The standard error of the estimate (48) of the mode interpreted as a grade, obtained by introducing  $y = 0$  into (49) is

(49') 
$$\sigma(u)\sqrt{n} = \frac{1}{\alpha} \sqrt{e - 1} = 1.02205\sigma.$$

The most precise grade is

$$\hat{x} = u + \frac{\hat{y}}{\alpha},$$

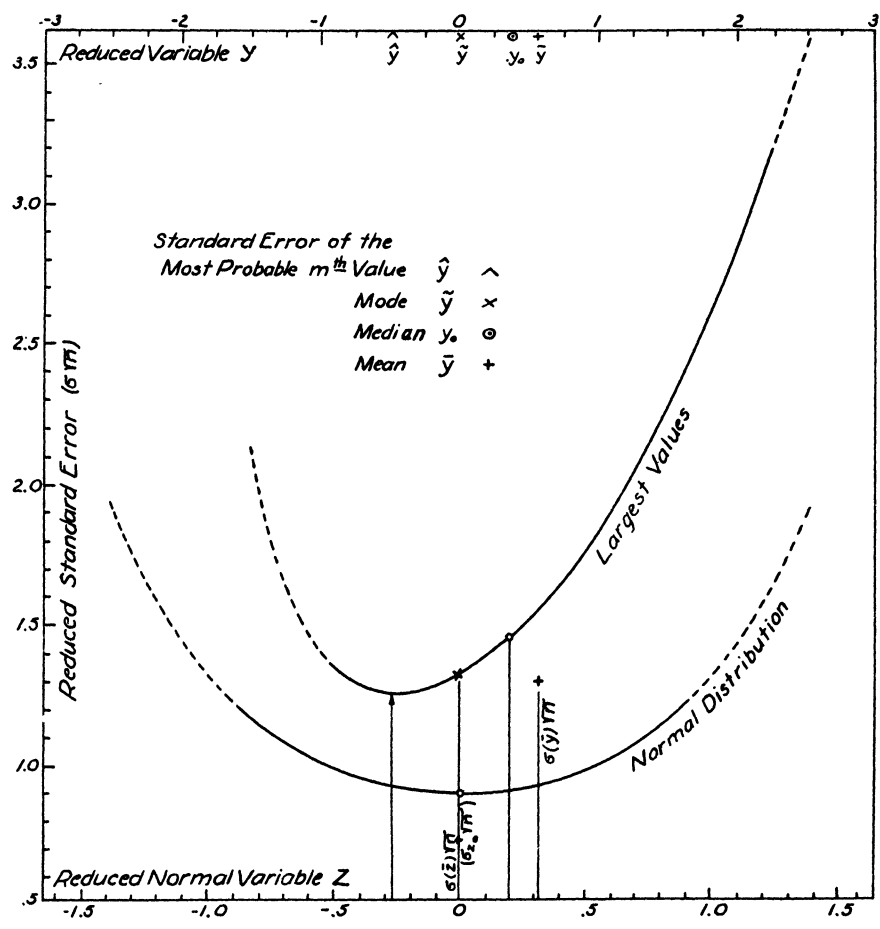


FIG. 1. The Reduced Standard Error of the  $m$ th Value

where  $\hat{y}$  is the value of the reduced variate, for which the standard error (49) is minimum. We obtain from (49) and (45) the numerical values

(50) 
$$\hat{y} = -.46601; G(\hat{y}) = .20319; \sigma(\hat{x})\sqrt{n} = .96887\sigma.$$

The standard error of the most precise grade is 3 per cent smaller; the standard error of the mode, estimated as a grade, is 2 per cent larger than the standard error of the mean.

**6. Confidence bands.** The standard errors (33) of the grades may be used in a general way for the construction of *confidence bands* obtained from curves which control the fit between theory and observation. Consider first the observed stepfunction  $(m - \frac{1}{2}, x_m)$  and the theoretical ogive  $nF(x)$ ,  $x$ . The variate  $x$  is plotted along the abscissa, the cumulative frequency along the ordinate. Now, for large  $n$  any theoretical value  $x$ , which is not extreme, may be interpreted as an  $m$ th value having a normal distribution and a standard error  $\sigma(x_m)$ . At each point of the graph of  $nF(x)$ ,  $x$  which is not extreme, we construct a segment of length  $2\sigma(x_m)$  parallel to the  $x$  axis, the midpoint of the segment being on the theoretical ogive. In other words, we add the standard error  $\sigma(x_m)$  to, and subtract it from, any corresponding value  $x$ , and attribute  $nF(x)$  to the beginning and end of these intervals. By this procedure we obtain two curves  $nF(x)$ ,  $x \mp \sigma(x_m)$ . For each observation there exists a probability  $P = .68268$  that it will be contained within the interval  $x \mp \sigma(x_m)$ .

If we apply another hypothesis to the same observations, or choose other values for the constants, we reach, of course, other control curves. Of two competing hypotheses the one is to be preferred where the band contains a larger number of observations.

The same method may be applied to the equiprobability test and to the comparison of the observed and theoretical return periods [6]. This procedure is legitimate for all values which are not extreme.

In the following, we construct the confidence bands for the normal distribution

$$(51) \quad g(z) = \frac{1}{\sqrt{\pi}} e^{-z^2}.$$

The variate  $x$  is related to the reduced variate  $z$  by (1), which, in this case, becomes

$$(52) \quad x = \bar{x} + \sigma\sqrt{2}z.$$

The probability  $G(z)$  is

$$(53) \quad G(z) = \frac{1}{2}(1 + \Phi(z)),$$

where  $\Phi(z)$  stands for the Gaussian integral

$$(54) \quad \Phi(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt.$$

Formulae (36) and (53) lead to the reduced standard error

$$(55) \quad \sigma(z)\sqrt{n} = \frac{1}{2g(z)} \sqrt{1 - \Phi^2(z)},$$

given in the table, col. 6. The standard errors  $\sigma(x_m)$  of the  $m$ th values obtained from (37) (52) and (55) are

$$(56) \quad \sigma(x_m) = \frac{\sigma\sqrt{2}}{\sqrt{n}} [\sigma(z)\sqrt{n}].$$

As a numerical example we choose the annual precipitations observed in 51 meteorological stations in Paris and its surroundings in the year 1938. We suppose that the differences between the 51 observations are only due to chance. The stepfunction  $m - \frac{1}{2}, x_m$  is plotted in figure 2. To obtain the theoretical ogive we compute the constants in (52). They are

$$(57) \quad \bar{x} = 571.92; \sigma\sqrt{2} = 38.52.$$

The theoretical values  $x$  obtained from (52), the cumulative frequencies  $nF(x)$  obtained from the table of the Gaussian integral [11] and the standard errors

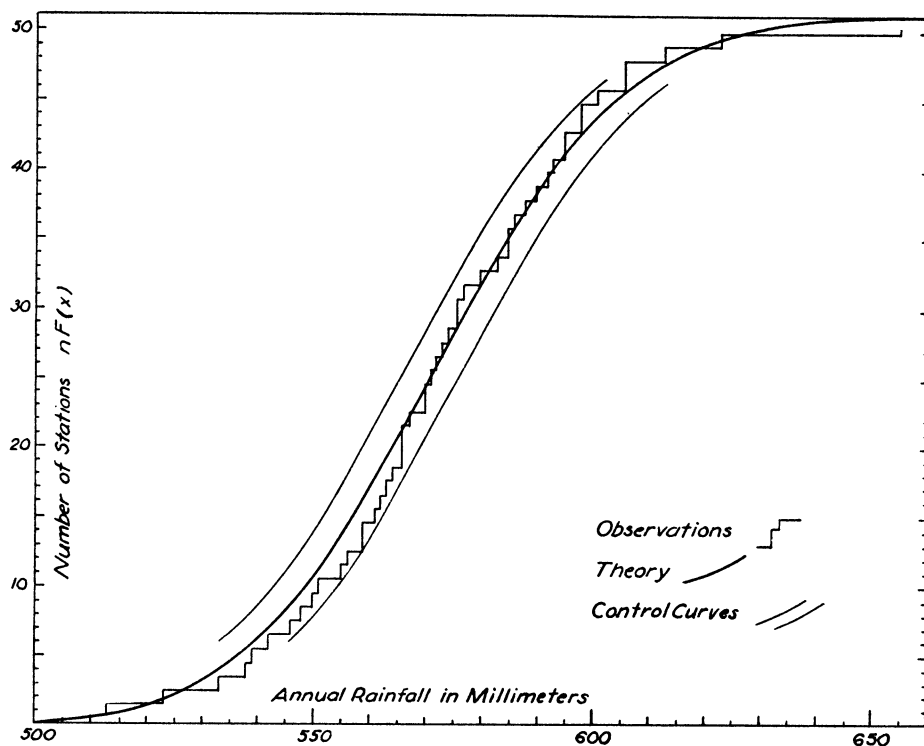


FIG. 2. The Confidence Band

$$(58) \quad \sigma(x_m) = 5.393 [\sigma(z)\sqrt{n}],$$

obtained from (56) are given in the columns 2 to 5 and 7 of Table I.

We trace in figure 2 the theoretical curve  $nF(x)$ ,  $x$  and the confidence band obtained from col. 7. by the methods described above. All observations are contained within the control curves. We may accept the theory that the differences between the annual rainfalls observed in the 51 stations are only due to chance.

**7. Conclusions.** To test a statistical hypothesis for a continuous variate we use the ogive, the equiprobability method, based on (1), and the return periods

(5). The three tests may be combined on appropriate probability paper. As the rank of the  $m$ th observation  $x_m$  may be  $m$  or  $m - 1$ , we have two series of observations. To obtain one and only one series we use for the ogive the serial number  $m - \frac{1}{2}$  provided that the number of observations is large. Generally, we attribute to  $x_m$  an adjusted frequency, namely, the probability (15) of the most probable  $m$ th value. The adjusted frequency is obtained from the serial number  $m - \frac{1}{2}$  and a correction,  $D$ , equation (17), which depends upon the distribution. The correction is important for the three tests, and small  $n$ , furthermore, for the equiprobability test and the return periods for the extreme observations and any number  $n$ .

The same correction  $D$  is used for estimating a grade through its relation (26) to the corresponding most probable serial number  $\tilde{m}$ . For distributions, where the second moment does not exist, we estimate the arithmetic mean from a

TABLE I  
*Normal Confidence Band and Theoretical Frequencies of the Rainfalls*

Reduced Variate $\pm z$ 1	Variate		Frequency		Reduced Standard Error $\sigma(z)\sqrt{n}$ 6	Standard Error $\sigma(x_m)$ 7
	$x_2$	$x_3$	$51 F(x)$ 4	$51 F(x)$ 5		
0	571.91	571.9	25.50	25.50	.886	4.8
.2	564.2	579.6	19.82	31.18	.899	4.9
.4	556.5	587.3	14.58	36.42	.940	5.1
.6	548.8	595.0	10.10	40.90	1.012	5.5
.8	541.0	602.7	6.58	44.42	1.127	6.1
1.0	533.4	610.4	4.01	46.99	1.297	7.0
1.2	525.7	618.1	2.29	48.71		
1.4	418.0	625.9	1.22	49.78		
1.6	510.3	633.6	.60	50.40		
1.8	502.6	641.3	.28	50.72		

grade. For asymmetrical distributions we estimate the mode from a grade by (32) and (48).

In this case, we have to introduce a distinction between the mode and the most precise grade (43). The adjusted frequency and the estimates for grades may be used even for small numbers of observations.

The standard error of these estimates is obtained, equation (33) from the limiting, normal, form of the distribution of the  $m$ th value, which holds, provided the serial number is not extreme. To control a hypothesis we construct confidence bands, which are obtained from the standard errors of the grades.

#### REFERENCES

- [1] R. A. FISHER AND L. H. C. TIPPETT, "Limiting forms of the frequency distribution of the largest or smallest member of a sample," *Proc. Camb. Phil. Soc.*, Vol. 24, part 2 (1928), p. 180.

- [2] E. J. GUMBEL, "Les valeurs extrêmes des distribution statistiques," *Annales de l'Institut Henri Poincaré*, Vol. 4 (1935), Paris, p. 115.
- [3] E. J. GUMBEL, "Les valeurs de position d'une variable aléatoire," *Comptes Rendus*, Vol. 208, (1939), Paris, p. 149.
- [4] E. J. GUMBEL, "The return period of flood flows," *Annals of Math. Stat.*, Vol. 12 (1942), p. 163.
- [5] E. J. GUMBEL, "Simple tests for given hypotheses," *Biometrika*, Vol. 32 (1942), p. 317.
- [6] E. J. GUMBEL, "Statistical control curves for flood discharge," *Trans. Am. Geoph. Union* (1942), Washington, p. 489.
- [7] ALLEN HAZEN, *Flood Flows*, New York, John Wiley, 1930.
- [8] B. F. KIMBALL, "Limited type of primary probability distribution applied to annual flood flows," *Annals of Math. Stat.*, Vol. 13 (1942), p. 318.
- [9] DIXON H. LEAVENS, "Frequency distributions corresponding to time series," *Jour. Amer. Stat. Assoc.*, Vol. 26 (1931), p. 407.
- [10] KARL AND MARGARET V. PEARSON, "On the mean character and variance of a ranked individual, and on the mean and variance of the interval between ranked individuals," *Biometrika*, Vol. 23, part 3, 4 (1931), p. 364; Vol. 24, part 1, 2 (1932), p. 203.
- [11] *Tables of Probability Functions*, Federal Works Agency W.P.A. of New York City, 1941.