

## 152

### THE WAVE OF ADVANCE OF ADVANTAGEOUS GENES

Author's Note (CMS 31.354a)

This is an isolated paper which I have not followed up either practically or theoretically. It seemed essential to examine the properties of the differential equation determining gene spread in the simplest case, and it was a pleasure to find that the very empirical process used for the computation of the fundamental function would work successfully. The quantitative solutions found were strongly suggestive for comparisons with the observable facts, especially in littoral organisms.

# THE WAVE OF ADVANCE OF ADVANTAGEOUS GENES

BY R. A. FISHER, Sc.D., F.R.S.

## I. THE PROBLEM OF GENE DISPERSION

CONSIDER a population distributed in a linear habitat, such as a shore line, which it occupies with uniform density. If at any point of the habitat a mutation occurs, which happens to be in some degree, however slight, advantageous to survival, in the totality of its effects, we may expect the mutant gene to increase at the expense of the allelomorph or allelomorphs previously occupying the same locus. This process will be first completed in the neighbourhood of the occurrence of the mutation, and later, as the advantageous gene is diffused into the surrounding population, in the adjacent portions of its range. Supposing the range to be long compared with the distances separating the sites of offspring from those of their parents, there will be, advancing from the origin, a wave of increase in the gene frequency. We may first on the simplest possible postulates consider the motion of this wave.

Let  $p$  be the frequency of the mutant gene, and  $q$  that of its parent allelomorph, which we shall suppose to be the only allelomorph present. Let  $m$  be the intensity of selection in favour of the mutant gene, supposed independent of  $p$ . Suppose that the rate of diffusion per generation across any boundary may be equated to

$$-k \frac{\partial p}{\partial x}$$

at that boundary,  $x$  being the co-ordinate measuring position in the linear habitat. Then  $p$  must satisfy the differential equation

$$\frac{\partial p}{\partial t} = k \frac{\partial^2 p}{\partial x^2} + mpq, \quad \dots\dots(1)$$

where  $t$  stands for time in generations.

The constant  $k$  is a coefficient of diffusion analogous to that used in physics. Its use should be appropriate in many cases. In all real cases we may expect irregularities due to  $k$  varying at different points of the range, due to variations in the density of the population, and to variation in the selective advantage of the mutant at different places. Further, the means of diffusion may involve an unequal drift in opposite directions, so that some parts of the range predominate as centres of multiplication and others as centres of extinction. The effects of all such complications can only be discussed by reference to the course of events when they are absent. The purpose of equation (1) is to specify the simplest possible conditions.

The use of the analogy of physical diffusion will only be satisfactory when the distances of dispersion in a single generation are small compared with the length of the wave. In reality diffusion is a complex process, compounded often of the diffusion of gametes, and that of

larvae, in addition to adult forms; a more exact treatment than that supplied by a simple coefficient would involve the interaction of these components, and the stages at which the selective advantage was enjoyed. So far as it is applicable, the analogy of physical diffusion, therefore, greatly simplifies the problem.

With respect to the assumed independence of  $m$  from  $p$ , this is effectively to assume that there is no dominance in respect of the selective advantage enjoyed. Apart from its simplicity this is also, in the author's opinion, the most important case to consider, in respect to advantageous mutations occurring in nature. There are, at least, plausible reasons for supposing that the common recessiveness of observed mutations is a characteristic of harmful mutations, which have long been appearing in the species with relatively high mutation rates, whereas beneficial mutations must, at the time of their establishment, occur with exceedingly low mutation rates, and have rarely appeared before in the recent history of the species. On these grounds dominance would be expected to be absent, and its absence is made the more probable by the fact that in most cases the quantitative effect of beneficial mutations must be extremely small. For the same reason the selective intensity  $m$  is taken to be a small quantity, so that  $p$  may be taken to vary continuously with time, and not discontinuously from generation to generation.

II. WAVES OF STATIONARY FORM

If we seek for a solution of (1) representing a wave of stationary form advancing with velocity  $v$ , we may put

$$\frac{\partial p}{\partial t} = -v \frac{\partial p}{\partial x},$$

and obtain the differential equation (2) involving only one independent variable:

$$k \frac{d^2 p}{dx^2} + v \frac{dp}{dx} + mpq = 0. \tag{2}$$

Since the variable  $x$  does not appear explicitly, we may write, for the frequency gradient,

$$g = -dp/dx,$$

whence

$$\frac{d^2 p}{dx^2} = -\frac{dg}{dx} = g \frac{dg}{dp},$$

and so find the relation between  $g$  and  $p$ ,

$$kg \frac{dg}{dp} - vg + mpq = 0. \tag{3}$$

At the point of inflexion  $dg/dp = 0$ , and  $vg = mpq$ ; in advance of this point  $dg/dp$  is positive. If  $g/p$  tends to a limit  $u$  as  $p$  tends to zero, then  $u$  must satisfy the equation,

$$ku^2 - vu + m = 0,$$

a quadratic equation in  $u$ , which has real roots only if  $v^2$  is not less than  $4km$ ; but  $g/p$  cannot tend to zero for  $vg > mpq$ , and cannot tend to infinity because  $v > k \cdot dg/dp$ . Hence solutions only exist for which the velocity of propagation is equal to, or exceeds,  $2\sqrt{(km)}$ .

Writing

$$\lambda = \sqrt{(k/m)},$$

$$v = \sqrt{(mk)} \left( c + \frac{1}{c} \right),$$

then equation (3) may be written

$$\lambda^2 g \frac{dg}{dp} - g\lambda \left( c + \frac{1}{c} \right) + pq = 0. \tag{4}$$

Or, if

$$\lambda g = pqz,$$

$$pq \frac{dz}{dp} + (1 - 2p)z - \left( c + \frac{1}{c} \right) + \frac{1}{z} = 0, \tag{5}$$

where  $c$  is any positive number, conventionally taken to be less than 1.

In the especially interesting case of minimal velocity, equation (5) may be written

$$pq \frac{dz}{dp} = 2pz - \frac{(1-z)^2}{z}; \tag{5a}$$

this case, when  $c = 1$ , we may call (a). If  $c$  lies between 1 and  $\sqrt{\frac{1}{2}}$ , we have a range of cases, which may be called case (b). When  $c = \sqrt{\frac{1}{2}}$  (case c), a second case of special interest arises with the equation

$$pq \frac{dz}{dp} = 2pz + \frac{3}{\sqrt{2}} - \left( z + \frac{1}{z} \right), \tag{5c}$$

and having a velocity of propagation  $\sqrt{\frac{9}{8}}$  times the minimum. Finally, in case (d),  $c$  is less than  $\sqrt{\frac{1}{2}}$ .

### III. PARTICULAR CASES

When  $p = 1$ , the only positive value of  $z$  for which  $dz/dp$  is finite is the positive root of the equation

$$z^2 + \left( c + \frac{1}{c} \right) z - 1 = 0$$

or 
$$z = \frac{1}{2c} \{ \sqrt{(c^4 + 6c^2 + 1)} - (c^2 + 1) \} = \alpha.$$

In the neighbourhood of all other values  $dz/dp$  increases inversely to  $(1-p)$ , so that no other finite value is admissible at  $p = 1$ . In general, at this extremity

$$(1-p)z \frac{dz}{dp} = z^2 + \left( c + \frac{1}{c} \right) z - 1$$

or 
$$\frac{dp}{1-p} = \frac{z dz}{z^2 + \left( c + \frac{1}{c} \right) z - 1} = \left( \frac{\alpha}{z-\alpha} - \frac{\beta}{z-\beta} \right) \frac{dz}{\alpha-\beta},$$

$$-\log q = \frac{c}{\sqrt{(c^4 + 6c^2 + 1)}} \log \frac{(z-\alpha)^\alpha}{(z-\beta)^\beta} + A,$$

which as  $-\log q \rightarrow \infty$  cannot be satisfied for any finite value of  $A$ .

The only admissible solution is therefore that for which  $z = \alpha$ , at the limit when  $p = 1$ .

When  $p = 0$ , we have in case (a)

$$p \frac{dz}{dp} = 2pz - \frac{1}{z} (1-z)^2.$$

For positive values of  $z$ , the right-hand side is positive when

$$(2p - 1)z^2 + 2z - 1$$

is positive; this is zero when  $z$  is  $z_1 = \frac{-1 \pm \sqrt{(2p)}}{2p - 1}$ .

When  $p > \frac{1}{2}$ , there is only one positive root. This root decreases (as  $p$  passes from 0 to 1) from 1 to  $(\sqrt{2} - 1)$ , which are the terminal values of  $z$ . Since  $dz/dp$  is positive when  $z > z_1$  (apart from a region of higher values when  $p < \frac{1}{2}$ ),  $z$  can never exceed  $z_1$  for intermediate values of  $p$ , for positive values of its derivative can never allow it to pass out of the region of positive values, so as to decrease to its final value. Consequently, in the neighbourhood of  $p = 0$ ,  $z$  must decrease even more rapidly than  $z_1$ . For small values of  $p$  therefore  $dz/dp$  must tend to a negative infinity. The differential equation to be satisfied in this region is

$$p \frac{dz}{dp} = -\frac{1}{z}(1-z)^2, \tag{6a}$$

or 
$$\frac{dp}{p} = \frac{-z dz}{(1-z)^2} = \frac{dz}{1-z} - \frac{dz}{(1-z)^2};$$

whence 
$$\log p = -\frac{1}{1-z} - \log(1-z) + A,$$

or 
$$p = \frac{A}{1-z} e^{-1/(1-z)}, \tag{7a}$$

where the constant of integration  $A$  is that which carries the solution to the terminal value  $z = \sqrt{2} - 1$ , at  $p = 1$ .

In case (b), since  $dz/dp$  is positive for all values of  $p$  when  $z$  lies between  $c$  and  $1/c$ , and since the terminal value  $\alpha$  is less than  $c$ , we must take  $z = c$  at  $p = 0$ ; we need a negative value of  $dz/dp$  at the terminus, satisfying the equation

$$p \frac{dz}{dp} = 2pc - \frac{1}{c} \left( \frac{1}{c} - c \right) (c - z), \tag{6b}$$

where  $c^2 > \frac{1}{2}$ . Writing the equation in the form

$$\frac{dz}{dp} - \frac{1-c^2}{pc^2} z = 2c - \frac{1-c^2}{cp},$$

or 
$$\frac{d}{dp} (zp^{1-1/c^2}) = 2cp^{1-1/c^2} - \frac{1-c^2}{c} p^{-1/c^2},$$

it appears that 
$$zp^{1-1/c^2} = \frac{2c}{2-1/c^2} p^{2-1/c^2} + cp^{1-1/c^2} - B,$$

in which again, the first term on the right-hand side is to be omitted, giving the solution

$$c - z = Bp^{1/c^2-1}. \tag{7b}$$

Since the power of  $p$  is less than unity,  $dz/dp$  is still infinite at the limit.

In the special case (c), where  $1/c^2 = 2$ , we find on integration

$$z/p = 2c \log p + c/p - c,$$

or

$$\sqrt{\frac{1}{2}} - z = p(c - \sqrt{2} \log p), \quad \dots\dots(7c)$$

tending to zero with  $p$ , but still with an infinite derivative.

Finally in case (d) 
$$c - z = \frac{2c^3}{1 - 2c^2} p.$$

In this case the constant of integration is associated with a negligible term. In fact by expanding  $c - z$  in powers of  $p$  as

$$z = c - \frac{2c^3}{1 - 2c^2} p + \beta p^2 + \gamma p^3 + \dots$$

and substituting, we have successive equations for  $\beta, \gamma, \dots$ , i.e.

$$\beta = \frac{2c^5(3 - 4c^2)}{(1 - 3c^2)(1 - 2c^2)^2},$$

$$\gamma = \frac{8c^7(5 - 11c^2 + 8c^4)}{(1 - 3c^2)(1 - 4c^2)(1 - 2c^2)^3}.$$

We have thus an expansion for  $z$  as a power series in  $p$ , a form of expansion which fails at the singular values  $c^{-2} = 3, 4, 5, \dots$ ; showing, nevertheless, that when  $c^{-2} > 2$ , the solution having  $z = c$  for  $p = 0$  is unique.

#### IV. THE AMBIGUITY OF VELOCITY

The most striking point about equation (2) is that the velocity of advance of the mutant factor appears to be indeterminate. If, for example, any part of the range were filled with the mutant form, and the zone of transition were artificially given frequencies with the low gradient of gene ratio appropriate to a high velocity, the mutation would spread with a higher velocity than if the initial gradient had been higher, and would continue to spread indefinitely with this higher velocity so long as uniform conditions were encountered. Common sense would, I think, lead us to believe that, though the velocity of advance might be temporarily enhanced by this method, yet ultimately, the velocity of advance would adjust itself so as to be the same irrespective of the initial conditions. If this is so, equation (2) must omit some essential element of the problem, and it is indeed clear that while a coefficient of diffusion may represent the biological conditions adequately in places where large numbers of individuals of both types are available, it cannot do so at the extreme front and back of the advancing wave, where the numbers of the mutant and the parent gene respectively are small, and where their distribution must be largely sporadic.

The effect of chance at the advancing front may be calculated by considering an aggregate of discrete particles, which increase in number with a relative growth rate  $m$ , as at the wave front of our original problem, but are free also to increase in numbers indefinitely in the interior of their range. We shall suppose them to be scattered at small unit intervals of time

ADVANTAGEOUS GENES

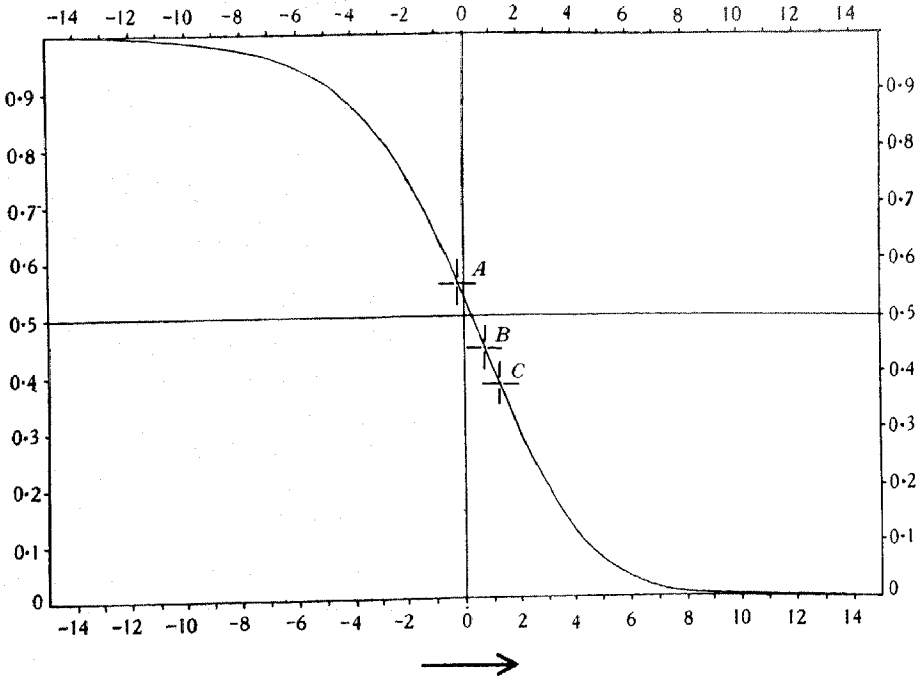


Fig. 1. Progressive wave of increase of frequency of advantageous genes. A, median of heterozygotes,  $p=0.559$ ,  $x=-0.194$ . B, point of inflexion, at which the rate of change of gene frequency is greatest,  $p=0.442$ ,  $x=0.765$ . C, point at which change in gene frequency is most easily detected,  $p=0.377$ ,  $x=1.297$ . The zero of the abscissa,  $x$ , is the point at which the number of mutant genes in front is equal to the number of parent genes behind,  $p=0.536$ ,  $x=0$ .

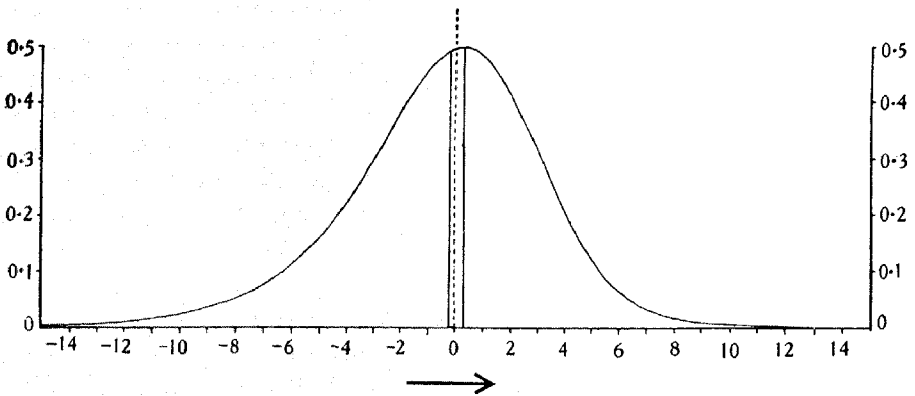


Fig. 2. Distribution of heterozygotes in relation to curve of increase of frequency of advantageous genes. Median  $x=-0.194$ ; mode  $x=+0.296$ .

so that the displacements of the particles at each scattering are distributed independently in the normal curve

$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\frac{x^2}{\sigma^2}} dx;$$

then  $k$  of our previous notation will correspond to  $\frac{1}{2}\sigma^2$ .

Whatever may be the original distribution, in one dimension, of the particles, we may specify it by means of the characteristic function

$$M(t) = S(e^{tx}),$$

where  $S$  stands for summation over all the particles, and  $x$  for the co-ordinate of any one of them. The effect of the dispersion of the particles is now merely to multiply  $M$  by the factor  $e^{\frac{1}{2}\sigma^2 t^2}$  at unit intervals of time, while the effect of multiplication of the particles is to multiply it by  $e^m$ . If  $K$  stands for  $\log M$ , it appears then that  $K$  increases uniformly with time at a rate  $m + \frac{1}{2}\sigma^2 t^2$ ; after time  $T$

$$K(T) = K(0) + T(m + \frac{1}{2}\sigma^2 t^2).$$

If the process be continued for a long time, the form of  $K$  will be determined by the ever-increasing second term, and the distribution will tend to the normal form, with variance  $T\sigma^2$ , and total number proportional to  $e^{mT}$ . Let us now draw a line beyond which a large but constant number of particles have already advanced, and consider with what velocity this line will move forward. The proportion,  $P$ , of the population beyond this line will fall off proportionately with  $e^{-mT}$ , but if  $\xi$  is the ratio of its distance from the centre to the standard deviation

$$P = \frac{1}{\xi\sqrt{2\pi}} e^{-\frac{1}{2}\xi^2}$$

approximately, when  $P$  is small, whence it appears that  $\xi^2$  differs from  $2mT$  by a constant, and by the logarithm of  $\xi$ , or, in other words, that  $\xi/\sqrt{2mT}$  tends to unity as a limit. But the ratio of the standard deviation to  $\sigma\sqrt{T}$  also tends to unity. Hence the distance of our arbitrary line from the centre bears a ratio to  $\sigma T\sqrt{2m}$ , which tends to unity. Evidently the front advances finally with constant velocity given by  $\sigma\sqrt{2m}$ , or putting  $\sigma = \sqrt{2k}$ , with velocity  $2\sqrt{km}$ , which is the minimal velocity consistent with equation (2). The conditions at the front of the wave are the same in both cases, save that the diffusion of a continuous variable has been replaced by the random dispersion of discrete particles, and when this is done it is seen that only one velocity of advance is ultimately possible.

### V. THE TABULATION OF THE WAVE FORM FOR $c = 1$

It has been shown in Section III that whereas innumerable solutions of the equation pass through the point  $p = 0, z = 1$ , only one passes through the other terminal point  $p = 1, z = \sqrt{2} - 1$ . Starting from this point therefore, it should be possible to obtain the numerical value of  $z$  for each value of  $p$  from 0 to 1, and so to construct the wave.

The process was carried out in three stages: (a) An expansion of  $z$  in terms of  $p$  was obtained



for the immediate neighbourhood of  $p = 1$ . (b) At any point on the curve  $dz/dp$  was calculated from the differential equation, and from this, and preceding values, the next point on the curve was obtained. (c) Since  $dz/dp$  tends to infinity as  $p$  tends to zero, at a certain stage a series of values of  $p$  for given  $z$  were obtained by interpolation, and the process continued using  $dp/dz$  instead of  $dz/dp$ .

$$\text{If } z = \sqrt{2} - 1 + aq + bq^2 + cq^3 + \dots$$

when  $q$  is small, by substitution in the differential equation, and equating powers of  $q$ , we find

$$a = \frac{2(\sqrt{2} - 1)}{5 + 2\sqrt{2}} = 0.10582293,$$

$$b = \frac{2(11\sqrt{2} - 1)}{(5 + 2\sqrt{2})^2(6 + 2\sqrt{2})} = 0.0533084,$$

$$c = \frac{16(35 + 38\sqrt{2})}{(5 + 2\sqrt{2})^3(6 + 2\sqrt{2})(7 + 2\sqrt{2})} = 0.0341074,$$

while the fourth coefficient is numerically about 0.02417. These suffice to give seven-figure accuracy up to  $q = 0.06$ , or numerically

Table I. *Values of  $z$  calculated from terminal expansion*

$p$	$z$
1.00	0.4142 136
0.99	0.4152 772
0.98	0.4163 518
0.97	0.4174 376
0.96	0.4185 3482
0.95	0.4196 4365
0.94	0.4207 6434

The first seven values for  $p$  and  $z$  give a sufficient start for the second process. From the value of  $z$  corresponding to  $p = 0.96$ , the value of  $dz/dp$  can be calculated from the differential equation

$$\frac{dz}{dp} = \frac{1}{pq} \left( \frac{1}{z} + 2qz - 2 - z \right).$$

Unit error in  $z$  will introduce an error in  $1/z$  of about 6, or in  $\frac{1}{z} - z$  of about 7; when the divisor  $pq$  is as small as 0.0384, the error in  $dz/dp$  is nearly 170 times as great as that in  $z$ , but in the opposite direction. As  $pq$  increases, however, the error in  $dz/dp$  becomes less than 100 times that in  $z$ , and the increment added to  $z$  to give the next value becomes sufficiently accurate.

The increment may be calculated from the differential coefficient, and its backward differences. This if  $D$  stand for the operation of differentiation,  $\Delta$  for forward differencing, and  $\nabla$  for backward differencing,

$$\Delta z = (e^D - 1)z = \left(1 + \frac{1}{2}D + \frac{1}{6}D^2 + \dots\right) Dz;$$

but

$$e^{-D} = 1 - \nabla,$$

$$D = \nabla + \frac{1}{2}\nabla^2 + \frac{1}{3}\nabla^3 + \dots;$$

whence

$$\Delta z = \left\{ 1 + \frac{1}{2}\nabla + \frac{5}{12}\nabla^2 + \frac{3}{8}\nabla^3 + \dots \right\} Dz;$$

which is conveniently applied in the form

$$1 + \frac{1}{2}\nabla \left( 1 + \frac{5}{6}\nabla \left( 1 + \frac{9}{10}\nabla \left( 1 + \dots \right) \right) \right)$$

where as many as three differences are used.

To minimize initial errors eight places were used in the values of  $z$  for  $p = 0.96, 0.95$  and  $0.94$ . The scheme of calculation then starts as below:

Table II. Calculation of  $z$  from the differential equation

$p$	$z$	$dz/dp$	$\nabla dz/dp$	$\nabla^2 dz/dp$
0.96	0.4185 3482	11029 8		
			117 4	
0.95	0.4196 4365	11147 2	119 9	2 5
0.94	0.4207 6434	11267 1	122 1	2 2
0.93	0.4218 9715	11389 2	124 8	2 7
0.92	0.4230 423	11514 0	127 8	3 0
0.91	0.4242 000	11641 8		
			129 4	1 6
0.90	0.4253 707	11771 2	133 2	3 8
0.89	0.4265 544	11904 4	135 5	2 3
0.88	0.4277 516	12039 9	138 6	3 1
0.87	0.4289 625	12178 5	141 6	3 0
0.86	0.4301 874	12320 1		

The second differences of  $dz/dp$  show but slight oscillation. It is not to be supposed that the seventh figure in  $z$  is always correct, but on trying a false start with an error  $-2$  at  $p = 0.96$ , and  $+4$  at  $p = 0.95$ , the errors in the subsequent figures alternate, with the greatest error  $-7$  at  $p = 0.93$ , and at  $p = 0.86$  are in exact agreement with the table above.

As the process is continued  $dz/dp$  and its differences increase. Third differences become appreciable at about  $p = 0.70$  and fourth differences at about  $p = 0.35$ . From  $p = 0.21$  to  $p = 0.15$  the interval was reduced to  $0.005$ , from  $0.15$  to  $0.10$  to  $0.002$ , and from  $0.10$  to  $0.07$  to  $0.001$ , in order to make the difference series decrease sufficiently rapidly.

From the values of  $z$  between  $p = 0.070$  and  $0.078$ , the values of  $p$  corresponding with  $z = 0.663, 0.664, 0.665$  and  $0.666$  were calculated, using initially nine figures, and from these the values of  $dp/dz$  calculated from the differential equation, thus

Table III. Values of  $p$  used in the final stages of tabulation

$z$	$p$	$dp/dz$		
0.663	0.0751002 75	0.968591	-9538	
0.664	0.0741364 63	0.959053		
			-9463	75
0.665	0.0731821 47	0.949590	-9387	76
0.666	0.0722372 61	0.940203	-9313	74
0.667	0.0713017 2	0.930890	-9240	73
0.668	0.0703754 6	0.921650	-9166	74
0.669	0.0694584 0	0.912484		

The values from  $z = 0.667$  onwards were obtained by calculating the successive differences from the differential coefficient. From  $z = 0.670$  to  $z = 0.790$  the interval was 0.002; but from that point using fourth differences the interval can be raised to 0.005. The last point calculated gave  $z = 0.875$ ,  $p = 0.000401738$ , at which stage  $p$  is decreasing by more than a quarter of its value at each step.

VI. NUMERICAL APPLICATIONS

It appears from equation (4) that the gradient is a maximum where  $g = pq/2\lambda$ , or where  $z = \frac{1}{2}$ . This occurs when  $p = 0.442428$ , when  $g\lambda = 0.1233427$ .

With a population cross-breeding at random the proportion of heterozygotes for any value of  $p$  is  $2pq$ . The total number of heterozygotes in any length of habitat in comparison with the number of organisms is

$$\int 2pq \, dx,$$

between the limits considered; writing  $\lambda dp/pqz$  for  $dx$ , this is seen to be merely

$$2\lambda \int \frac{dp}{z}.$$

Now

$$d(pqz) = pq \, dz + (1 - 2p)z \, dp,$$

but by equation (5 a)

$$pq \, dz = 2pz \, dp - \frac{dp}{z}(1 - z)^2;$$

hence

$$\begin{aligned} d(pqz) &= \frac{dp}{z} \{z^2 - (1 - z)^2\} \\ &= 2dp - \frac{dp}{z}. \end{aligned}$$

Consequently, the indefinite integral

$$2\lambda \int \frac{dp}{z}$$

may be expressed in the form

$$2\lambda(2p - pqz).$$

Since  $pqz$  vanishes when  $p = 0$ , or  $p = 1$ , the total number of heterozygotes maintained at any time is equal to the population of the length of habitat

$$4\lambda = 4\sqrt{(k/m)}$$

proportional to the square root of the coefficient of diffusion, and inversely to the square root of the intensity of selection. The relation

$$\int \frac{dp}{z} = 2p - pqz$$

also affords a needed check to the accuracy of the values of  $z$  obtained, for if the process of calculation at any stage had allowed of a systematic drift across the curves satisfying the equation, the values of  $1/z$  would have been systematically raised or lowered, and the value of the integral would depart from its calculated value. The test may be applied by sections.

The largest discrepancy found is between 0.4 and 0.5, and amounts to about one part in two millions, or nearly  $2\frac{1}{2}$  units in the seventh place. The value of  $p$  at the point of inflexion may thus really be one part in a million higher than that given above.

The effective centre of the wave in its advance is the point at which there are as many mutant genes in front as there are parent genes behind. The number of parent genes behind any point, expressed in terms of the population per unit length of the habitat, is

$$\int_{-\infty}^1 q dx = \int^1 \frac{q dp}{g} = \lambda \int^1 \frac{dp}{pz}$$

Since  $z$  behaves regularly in the neighbourhood  $p = 1$ , this integral offers no difficulty to direct evaluation. The value from  $p = \frac{1}{2}$  comes to  $1.540762\lambda$ . The number of mutant genes in advance of any point is more troublesome to ascertain. The form

$$\lambda \int_0^1 \frac{dp}{qz}$$

is unsuitable, since the differential coefficient of  $z$  is infinite at  $p = 0$ . Writing

$$2 - \frac{d}{dp}(pqz)$$

for  $1/z$ , it takes the form  $\lambda \left( -2 \log q - \int_0^1 \frac{1}{q} d(pqz) \right)$ ,

which may be used, though with some difficulty. Near the terminus, the most satisfactory process is to expand  $d(pqz)$  in the form

$$d(pqz) = pz dq + qd(pz),$$

giving the third form

$$\lambda \left( -2 \log q - pz - \int pz \frac{dq}{q} \right) = \lambda \left( -2 \log q - pz - (p + \log q)z - \int^1 (p + \log q) dz \right),$$

which may be used with confidence, since  $p + \log q$  is of the order of  $\frac{1}{2}p^2$ , and becomes negligible within the range tabulated.

The integral to  $p = \frac{1}{2}$  is found to be  $1.244939\lambda$ , showing, on comparison with the number of parent genes behind this point, that the effective centre lies behind the 50 per cent point by  $0.295823\lambda$ , or at a place where  $p$  has risen to  $0.535709$ . At this point the number of mutant genes in front and of parent genes behind are both equal to the total number in the length  $1.398150\lambda$ . We take this point as the origin of the co-ordinate  $x$ , in Table IV.

To put the situation concretely, let us suppose a mutation giving a selective advantage of 1 per cent is spreading along a continuously occupied shore line. Suppose that the standard displacement of young from parents in each generation is 100 yards. Then with  $m = 0.01$  per generation,  $k = 5000$  square yards per generation, and  $\lambda = \sqrt{(k/m)} = 717$  yards. The number of heterozygotes is equal to the population of 2868 yards, or rather more than a mile and a half of coast, though it is spread over 6 or 8 miles. The rate of advance  $v = 2\sqrt{(mk)}$  is

Table IV. Values of  $z$  and  $x$  for each integral percentage of  $p$ .  
 The gradient,  $-dp/dx$ , is  $pqz$ , when  $q = 1 - p$

$p$	$z$	$x$	$p$	$z$	$x$
0	1.000 0000	$\infty$			
0.01	0.784 0738	7.6061	0.50	0.487 3135	0.2958
0.02	0.749 9772	6.6817	0.51	0.485 2501	0.2136
0.03	0.726 5592	6.1182	0.52	0.483 2244	0.1309
0.04	0.708 2499	5.7027	0.53	0.481 2350	0.0477
0.05	0.693 0407	5.3693	0.54	0.479 2807	- 0.0360
0.06	0.679 9482	5.0883	0.55	0.477 3604	- 0.1203
0.07	0.668 4082	4.8437	0.56	0.475 4729	- 0.2053
0.08	0.658 0635	4.6261	0.57	0.473 6171	- 0.2910
0.09	0.648 6717	4.4291	0.58	0.471 7921	- 0.3776
0.10	0.640 0603	4.2485	0.59	0.469 9969	- 0.4651
0.11	0.632 1012	4.0811	0.60	0.468 2305	- 0.5536
0.12	0.624 6967	3.9246	0.61	0.466 4921	- 0.6431
0.13	0.617 7705	3.7774	0.62	0.464 7808	- 0.7338
0.14	0.611 2612	3.6380	0.63	0.463 0959	- 0.8358
0.15	0.605 1194	3.5053	0.64	0.461 4365	- 0.9191
0.16	0.599 3038	3.3785	0.65	0.459 8020	- 1.0139
0.17	0.593 7802	3.2568	0.66	0.458 1916	- 1.1103
0.18	0.588 5197	3.1396	0.67	0.456 6047	- 1.2085
0.19	0.583 4974	3.0264	0.68	0.455 0406	- 1.3085
0.20	0.578 6922	2.9167	0.69	0.453 4987	- 1.4105
0.21	0.574 0856	2.8102	0.70	0.451 9785	- 1.5147
0.22	0.569 6611	2.7066	0.71	0.450 4792	- 1.6213
0.23	0.565 4051	2.6056	0.72	0.449 0005	- 1.7304
0.24	0.561 3047	2.5068	0.73	0.447 5417	- 1.8423
0.25	0.557 3488	2.4101	0.74	0.446 1024	- 1.9572
0.26	0.553 5274	2.3154	0.75	0.444 6820	- 2.0754
0.27	0.549 8316	2.2223	0.76	0.443 2802	- 2.1972
0.28	0.546 2533	2.1308	0.77	0.441 8964	- 2.3229
0.29	0.542 7852	2.0406	0.78	0.440 5303	- 2.4529
0.30	0.539 4208	1.9518	0.79	0.439 1812	- 2.5876
0.31	0.536 1539	1.8640	0.80	0.437 8492	- 2.7276
0.32	0.532 9792	1.7773	0.81	0.436 5330	- 2.8733
0.33	0.529 8914	1.6915	0.82	0.435 2329	- 3.0255
0.34	0.526 8861	1.6066	0.83	0.433 9492	- 3.1849
0.35	0.523 9589	1.5224	0.84	0.432 6805	- 3.3525
0.36	0.521 1060	1.4388	0.85	0.431 4265	- 3.5293
0.37	0.518 3237	1.3558	0.86	0.430 1874	- 3.7166
0.38	0.515 6085	1.2732	0.87	0.428 9625	- 3.9160
0.39	0.512 9575	1.1911	0.88	0.427 7516	- 4.1295
0.40	0.510 3676	1.1093	0.89	0.426 5544	- 4.3597
0.41	0.507 8362	1.0278	0.90	0.425 3706	- 4.6097
0.42	0.505 3607	0.9465	0.91	0.424 2001	- 4.8837
0.43	0.502 9387	0.8653	0.92	0.423 0422	- 5.1876
0.44	0.500 5680	0.7842	0.93	0.421 8972	- 5.5293
0.45	0.498 2466	0.7031	0.94	0.420 7643	- 5.9205
0.46	0.495 9724	0.6220	0.95	0.419 6436	- 6.3796
0.47	0.493 7437	0.5408	0.96	0.418 5348	- 6.9371
0.48	0.491 5586	0.4594	0.97	0.417 4376	- 7.6502
0.49	0.489 4159	0.3777	0.98	0.416 3518	- 8.6479
0.50	0.487 3135	0.2958	0.99	0.415 2772	- 10.3480
			1.00	0.414 2136	$\infty$

about 14 yards per generation, or less than 10 miles in 1000 generations. To spread over a habitat of several hundred miles might well take 10,000 or 100,000 generations. In consequence, at any one time, the number of such waves of selective advance, simultaneously in progress, must be large. The effective centre in our example is about 210 yards behind the 50 per cent point, while the steepest gradient of gene ratio, which is the point of most rapid genetic change, is about 330 yards in advance of this point.

At any given spot the rate of change per generation in the proportion of mutant genes is

$$vg = 2pqmz,$$

which is less than  $\frac{1}{4}$  per cent at its highest point, where  $p$  is about 44 per cent. Very large counts would therefore be needed, supposing the gene to affect any measurable or observable characteristics, to detect the change in progress by observations during the course of only a few generations. If, for example, both homozygotes and heterozygotes could be distinguished with certainty, the sampling variance of  $p$ , as estimated from the examination of  $n$  individuals, would be  $pq/2n$ , while that of the difference as estimated from two such counts would be  $pq/n$ .

If  $n$  were as high as 10,000, the standard error is thus  $\frac{1}{2}$  per cent when  $p$  is 0.5 where the rate of change is only 0.244 per cent in each generation, so that about 5 generations must elapse before a significant increase in the percentage could be observed. The rate of change is greatest in relation to its sampling error at the maximum value of  $z\sqrt{pq}$  or when

$$z = \frac{1}{1 + \sqrt{\frac{1}{2} + p}},$$

which occurs when  $p = 0.377$ , or about  $1.297\lambda$  in advance of the effective centre of the wave, where, if the number counted,  $n$ , is equated to  $1/m^2$ , the rate of change is just over half (0.5005) its standard deviation in each generation. If the change manifests itself in a metrical character, to the variance of which other factors, environmental or genetic, contribute, the change will be most easily detected at some point between  $P = 0.377$  and  $P = 0.442$ . The direction of advance might also be indicated from observations at a single epoch by the asymmetry of the wave, which is more extended behind than before, or by the skewness of the distribution of heterozygotes, though these features might be expected to be obscured by irregularities in the habitat.

## VII. APPENDIX ON THE CALCULATION OF SPECIAL POINTS

### (a) *The point of inflexion*

At the point of greatest gradient, since in general

$$kg \frac{dg}{dp} - vg + mpq = 0,$$

we have the relation

$$vg = mpq,$$

or, simply,

$$z = \frac{1}{2}.$$

The relevant tabular values are

$p$	$z$	$\delta^2 z$	$x$	$\delta^2 x$
0.46	0.4959724	—	0.622002	—
0.45	0.4982466	472	0.703128	-47
0.44	0.5005680	493	0.784207	+16
0.43	0.5029387	—	0.865302	—

Inverse interpolation for  $z = 0.5$  gives  $p = 0.4424276$ ; whence direct interpolation gives  $x = +0.764525$ .

(b) *The point at which changes of frequency are most easily detected*

This will not be at the point where change in frequency is most rapid, because the standard error of a comparison of frequencies is not constant, but varies as  $\sqrt{(pq)}$ . We must therefore maximize not  $zpq$ , but  $z\sqrt{(pq)}$ . This gives

$$\frac{dz}{dp} = -\frac{1}{2}z\left(\frac{1}{p} - \frac{1}{q}\right);$$

but, in general,

$$pq \frac{dz}{dp} = 2 - \frac{1}{z} - (1 - 2p)z,$$

hence

$$(1 - 2p)z^2 - 4z + 2 = 0,$$

or

$$z = \frac{1}{1 + \sqrt{(\frac{1}{2} + p)}}.$$

The numerical values in this neighbourhood are

$p$	$\frac{1}{1 + \sqrt{(\frac{1}{2} + p)}}$	$z$	Difference	$\delta^2$
0.36	0.5188439	0.5211060	-0.0022621	—
0.37	0.5174007	0.5183237	-0.0009230	-509
0.38	0.5159737	0.5156085	+0.0003652	-471
0.39	0.5145638	0.5129575	+0.0016063	—

Inverse interpolation for zero difference gives  $p = 0.377126$ ; whence  $x$  is found to be 1.297092.

(c) *The median of the heterozygotes*

Since the proportion of the heterozygotes behind any point is given by

$$p - \frac{1}{2}pqz,$$

we may equate this expression to  $\frac{1}{2}$ . The following values will serve for interpolation:

$p$	$p - \frac{1}{2}pqz$	$\delta^2$	$x$	$\delta^2$
0.54	0.48047333	—	-0.035989	—
0.55	0.49092665	4178	-0.120301	-680
0.56	0.50142174	4139	-0.205293	-752
0.57	0.51195822	—	-0.291037	—

Inverse interpolation gives  $p = 0.5586476$ ,  $x = -0.193757$ .

## VIII. SUMMARY

The form is discussed of a steadily progressive wave of gene increase due to the local establishment of a favourable mutation, for the case of a uniform linearly distributed population.

The equation obtained by the analogy of physical diffusion is found to be consistent with all velocities of advance above a certain lower limit.

The indeterminacy of velocity is resolved by comparison with the properties of multiplying aggregates of particles, constantly subjected to random scattering. It appears that the actual velocity of advance must be the minimum compatible with the differential equation.

This velocity is proportional to the square root of the intensity of selective advantage and to the standard deviation of scattering in each generation, or to the square root of the diffusion coefficient when time is measured in generations. It may be expressed in the form

$$v = \sigma \sqrt{2m},$$

or

$$v = 2\sqrt{km},$$

where  $m$  is the selective advantage,  $\sigma$  the standard deviation of scattering, and  $k$  the diffusion coefficient.

The "length" of the wave, or the distance between any two assigned gene ratios, is proportional to

$$\lambda = \sqrt{k/m},$$

which may conveniently be taken as the unit of length.

The form of the wave is tabulated so as to show, for each percentage of the frequency of the mutant gene, the value of the gradient of gene ratio and the position at which this percentage occurs relative to the effective centre of the wave, i.e. to the point in advance of which there are as many mutant genes as there are parent genes behind it.

Stages of special interest which occur in succession at each point reached are

	Mutant genes %	Distance in advance of centre
The point at which changes in frequency are most easily detected	37.7	1.30 $\lambda$
The point of inflexion	44.2	0.76 $\lambda$
Equality of gene ratio, mode of heterozygotes	50.0	0.30 $\lambda$
Effective centre of wave	53.6	0
Median of heterozygotes	55.9	-0.19 $\lambda$