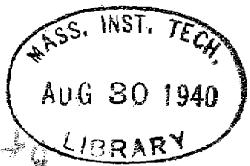


*ANALYSIS*  
*ANALYSIS*  
*BY CLAUDE ELWOOD SHANNON*

AN ALGEBRA FOR THEORETICAL GENETICS



By  
Claude Elwood Shannon  
B.S., University of Michigan  
1936

Submitted in Partial Fulfillment of the  
Requirements for the Degree of  
Doctor of Philosophy

From The  
Massachusetts Institute of Technology  
1940

Signature of Author.....  
Department of Mathematics, April 15, 1940  
Signature of Professor  
in Charge of Research.  
Signature of Chairman of Department  
Committee on Graduate Students....  
/

**MIT Document Services**

Room 14-0551  
77 Massachusetts Avenue  
Cambridge, MA 02139  
ph: 617/253-5668 | fx: 617/253-1690  
email: [docs@mit.edu](mailto:docs@mit.edu)  
<http://libraries.mit.edu/docs>

**DISCLAIMER OF QUALITY**

Due to the condition of the original material, there are unavoidable flaws in this reproduction. We have made every effort to provide you with the best copy available. If you are dissatisfied with this product and find it unusable, please contact Document Services as soon as possible.

Thank you.

## Table of Contents

I.	Introduction.....	1
II.	Notation.....	9
III.	Fundamental Theorems.....	16
IV.	The Solution of Equations Involving Unknown Populations.....	48
V.	Lethal Factors and Selection.....	53
VI.	A Calculus of Populations.....	59
	Bibliography.....	63
	Biography of the Author.....	64

238484

## I. Introduction

In this paper an attempt will be made to develop an algebra especially suited to problems in the dynamics of Mendelian populations. Many of the results presented here are old in the theory of genetics, but are included because the method of proof is novel, and usually simpler and more general than those used previously.

For the benefit of readers who are not familiar with modern genetics theory, we will first give a brief factual summary of those parts of it which are necessary for our work. Although all parts of the theory have not been uncontestedly established, still it is possible for our purposes, to act as though they were, since the results obtained are known to be the same as if the simple representation which we give were true. Hereafter we shall speak therefore as though the genes actually exist and as though our simple representation of hereditary phenomena were really true, since so far as we are concerned, this might just as well be so. We will omit from consideration mutations and phenomena in the sex chromosomes.

Heredity traits are transmitted by small elements called genes. These genes are carried in rodlike bodies known as chromosomes, a large number of genes lying side by side along the length of a chromosome. Chromosomes occur in pairs and an individual obtains one chromosome of each pair from his mother and the other from his fa-

ther.

By the genetic constitution of an individual we mean the kind and location of the genes which he possesses. If we represent the different genes by letters, then we may write a genetic formula for an individual. Thus considering two chromosome pairs and four gene positions in each chromosome, an individual might have the formula:

$$\begin{array}{ll} A_1 B_1 C_3 D_5 & E_4 F_1 G_6 H_1 \\ A_3 B_1 C_4 D_3 & E_4 F_2 G_6 H_2 \end{array} \quad (1)$$

Here the series  $A_1 B_1 C_3 D_5$  represents one chromosome, with  $A_3 B_1 C_4 D_3$  the corresponding one of the first pair.  $A_1, B_1, C_3, D_5, A_3, B_1, C_4, D_3$  are the genes lying in the positions under consideration.  $E_4 F_1 G_6 H_1$  and  $E_4 F_2 G_6 H_2$  are the two chromosomes of the second pair. We will sometimes write a genetic formula in one line. Thus (1) would be written:

$$A_1 A_3 B_1 B_1 C_3 C_4 D_5 D_3 E_4 E_4 F_1 F_2 G_6 G_6 H_1 H_2$$

alternate letters being taken from the top and bottom lines of (1).

There is no essential ordering of chromosomes in a pair. That is to say that the top and bottom lines of the formula for a chromosome pair may be inverted and still represent the same individual. Thus the formula (1) is identical, for example with the following:

A <sub>3</sub>	B <sub>1</sub>	C <sub>4</sub>	D <sub>3</sub>	E <sub>4</sub>	F <sub>1</sub>	G <sub>6</sub>	H <sub>1</sub>
A <sub>1</sub>	B <sub>1</sub>	C <sub>3</sub>	D <sub>5</sub>	E <sub>4</sub>	F <sub>2</sub>	G <sub>6</sub>	H <sub>2</sub>

in which we have inverted the first pair.

Certain simple traits are controlled by only one pair of genes lying at analogous points in corresponding chromosomes. Two such corresponding points in a chromosome pair are known as a gene locus, and the different genes which may occupy one locus are known as allelomorphs or more shortly as alleles. In our example (1) the positions occupied by genes C<sub>2</sub> and C<sub>4</sub> constitute a locus. We shall adopt the convention that allelomorphic genes shall have the same base letter with different subscripts. Thus C<sub>1</sub>, C<sub>2</sub>, C<sub>3</sub>, C<sub>4</sub>, C<sub>5</sub> represent five alleles. A C gene can only occur in the locus corresponding to C genes.

The appearance of an individual depends only on the kinds of genes, not on their positions. Thus an individual with the formula

A <sub>1</sub>	B <sub>1</sub>	C <sub>4</sub>	D <sub>3</sub>	E <sub>4</sub>	F <sub>1</sub>	G <sub>6</sub>	H <sub>2</sub>
A <sub>3</sub>	B <sub>1</sub>	C <sub>3</sub>	D <sub>5</sub>	E <sub>4</sub>	F <sub>2</sub>	G <sub>6</sub>	H <sub>1</sub>

would appear (insofar as the characteristics controlled by these genes are concerned) the same as (1). He would, however, breed differently as will appear later. Two such individuals are said to be phenotypically the same with respect to these characteristics. They are genotypically different; they have different genetic formulae with respect

to these loci. Such a situation can occur in a different way. In garden peas there are two alleles which control the size of the plant. These genes we may represent by  $A_1$  and  $A_2$ . If a plant has two  $A_1$  genes, it will be tall. If it has two  $A_2$  genes, it will be a dwarf. A plant with one  $A_1$  gene and one  $A_2$  gene is tall, since the gene for tallness ( $A_1$ ) is, as we say, dominant over the recessive gene ( $A_2$ ) for shortness. Thus,  $A_1A_1$  plants and  $A_1A_2$  (or  $A_2A_1$ ) plants are phenotypically the same but genotypically different with respect to tallness.

As we stated above, an individual receives one chromosome of each pair from the corresponding pair possessed by his mother and the other from that of his father. Let us now consider a pair of chromosomes possessed by a parent. In case a phenomenon known as cross-over does not occur in the chromosome pair under consideration, an offspring receives an entire chromosome selected at random from these two. We say that the genes in the chromosome are linked together meaning that they tend to be transmitted as a body. Genes located close together in the same chromosome are closely linked; the greater the distance between them, the weaker the linkage. Let us suppose that an individual has the genetic formula represented as follows

$$\begin{array}{ccccccc} A_1 & B_2 & C_3 & D_3 & E_4 & F_1 & G_1 \\ A_2 & B_2 & C_2 & D_1 & E_6 & F_1 & G_2 \end{array}$$

for a pair of corresponding chromosomes. Now, as we have said, in case crossover does not occur, an offspring of this individual will receive either the series

$$A_1 B_2 C_3 D_3 E_4 F_1 G_1 \text{ or } A_2 B_2 C_2 D_1 E_6 F_1 G_2$$

and he is equally likely to receive either of these.

However, it may happen that a crossover occurs between these chromosomes. If this crossover occurred, for instance, between the C and D loci, he would receive either

$$A_1 B_2 C_3 D_1 E_6 F_1 G_2 \text{ or } A_2 B_2 C_2 D_3 E_4 F_1 G_1$$

There is a definite probability that a crossover will occur between any two gene loci. Determining the relative positions of genes in a chromosome according to such a probability scale is known as mapping the chromosome. This has been carried out quite extensively for *Drosophila* and to a lesser extent for some other plants and animals. The map distance between two loci a and b may be defined as follows. Let  $x$  measure the actual physical distance along the chromosome and let  $p(x)$  be the probability that a crossover occurs between the points  $x$  and  $x + dx$ , providing it is known that no other crossover occurs near to the point  $x$ . This last restriction is necessary due to a phenomenon known as interference in which a crossover at one point hinders nearby crossovers. The map distance is then given by

$$\int_a^b p(x) dx.$$

The recombination value of two loci is the probability of an odd number of crossovers between these loci. For small distances the probability of more than one crossover is a second order term and the map distance is nearly equal to the recombination value, and both approximate the probability of one crossover between the loci.

If the two genes in a certain locus are identical, the individual is said to be homozygous in this factor. Otherwise he is heterozygous. The individual (1) is thus homozygous in the B, E, and G factors and heterozygous in all others.

A simple example will perhaps help to clarify these notions. Suppose that two gene loci are under consideration. There are three allelomorphic genes for the first locus,  $A_1, A_2, A_3$ ; the second locus has four alleles,  $B_1, B_2, B_3, B_4$ . The recombination value for these two loci is  $1/4$ . An individual with the genetic formula:

$$\begin{matrix} A_1 & B_4 \\ A_3 & B_2 \end{matrix}$$

is mated with an individual having the formula:

$$\begin{matrix} A_{\cancel{3}} & B_4 \\ A_2 & B_1 \end{matrix}$$

What is the probability that an offspring of this mating will have the formula

$$\begin{matrix} A_3 & B_4 \\ A_2 & B_1 \end{matrix}$$

Stated another way, what fraction of the offspring population should be expected to have this formula?

Evidently an offspring must obtain the  $A_3 B_4$  chromosome from the first parent. The probability that he will get an  $A_3$  gene from this parent is  $1/2$  since  $A_3$  and  $A_1$  are equally likely. If he gets this  $A_3$  gene, the probability that he will also get a  $B_4$  gene from this parent is  $1/4$ , the recombination value, since  $A_3$  and  $B_4$  are in opposite chromosomes. Thus, the probability that both events occur is  $1/2 \cdot 1/4 = 1/8$ . Now our offspring must obtain  $A_2$  and  $B_1$  from the second parent. He will certainly obtain an  $A_2$  since both genes in this locus are of this type. The chance that he obtains a  $B_1$  is  $1/2$ , since  $B_1$  and  $B_3$  are equally likely. The probability of the combination is therefore also  $1/2$ . Our final answer is, since the events are independent,  $1/8 \cdot 1/2 = 1/16$ . If we had asked what fraction would be of the type

$$\begin{matrix} A_3 & B_2 \\ A_2 & B_2 \end{matrix}$$

then in place of multiplying by the recombination value

1/4, we would multiply by  $1 - \frac{1}{4} = \frac{3}{4}$  since this is the probability that a crossover does not occur between the loci.

## II. Notation

To non-mathematicians we point out that it is a commonplace of modern algebra for the symbols to represent concepts other than numbers, and frequently therefore not to obey all the laws governing numbers. Such is the case in vector algebra, the theory of groups, rings, matrix algebra, in symbolic logic, tensor analysis, etc. In the particular algebra we construct for genetics theory the symbols represent Mendelian populations, and stand for a large group of numbers which describe the genetic constitution of the population. Addition and multiplication are defined to mean simple combination and cross breeding respectively, and it is shown that nearly all the laws of ordinary numerical algebra hold here. One interesting exception is the associative law of multiplication. It is not in general true that

$$(\lambda \times \mu) \times \nu = \lambda \times (\mu \times \nu)$$

Much of the power and elegance of any mathematical theory depends on use of a suitably compact and suggestive notation, which nevertheless completely describes the concepts involved. We will employ an index notation somewhat similar to that of the tensor calculus, which has proven so useful in differential geometry and in relativity.

theory. Because the notation employed is so basic to our work we will first explain the meaning of indexed symbols.

Consider, for example, the symbol

$$\lambda_{j k}^{h i} \quad (2)$$

Here  $\lambda$  is the base letter and  $h$ ,  $i$ ,  $j$ , and  $k$  are indices. Each index has a certain specific range of variation and the different indices may vary independently and even have different ranges of variation. In our work two indices in the same vertical column, such as  $h$  and  $j$  in (2), will always have the same range, but vary independently over this range. Thus  $h$  and  $j$  might have the range of values 1, 2, and 3 while  $i$  and  $k$  have the range 1, 2, 3, 4...9.

When the indices of (1) take on specific values, e.g.  $h = 1$ ,  $j = 3$ ,  $i = 5$ ,  $k = 5$ , the symbol

$$\lambda_{35}^{15}$$

represents a number. Symbol (2) then stands for a whole group of numbers, one for each combination of values of the indices; however, it should not be thought of as a group of separate numbers, but rather as a single entity having components whose values are the different numbers of the array.

When we think of an indexed symbol as representing a whole array of numbers and the indices as variables which assume any of the values in their ranges we say the indices are live or variable. Occasionally, however, it is desirable to think of  $\lambda_{j,k}^{hi}$  (say) as representing a certain specific one of the components. Thus we may set  $h=1, i=3, j=2, k=3$ . We say then that we have fixed or killed the indices at these values; they become dead indices. Also we sometimes wish to think of the indices as fixed at some value which is perfectly arbitrary. Without any change of notation we use  $\lambda_{j,k}^{hi}$  to represent an arbitrary component rather than the whole set of components. In such a case fixing the indices is purely subjective.

In an equation, although indices represented by different letters may vary independently, a specific letter e.g.  $h$ , must not take on different values in different places. Thus the sum of two indexed symbols

$$\lambda_{j,k}^{hi} + \mu_{j,k}^{hi} \quad (3)$$

is an indexed symbol, say  $\nu_{j,k}^{hi}$ , whose components are the sums of the corresponding components of  $\lambda_{j,k}^{hi}$  and  $\mu_{j,k}^{hi}$ . For example

$$\nu_{1,1}^{11} = \lambda_{1,1}^{11} + \mu_{1,1}^{11}$$

$$\nu_{3,5}^{12} = \lambda_{3,5}^{12} + \mu_{3,5}^{12}.$$

On the other hand, if

$$v_{j\kappa}^{ki} = \lambda_{jk}^{ki} + \mu_{jk}^{ki}$$

$$\text{then } v_{35}^{12} = \lambda_{35}^{12} + \mu_{35}^{12} \text{ etc.}$$

An equation in indexed symbols stands therefore for a large number of ordinary equations, one for each combination of values of the variable indices.

Ordinary multiplication of indexed symbols will be indicated by juxtaposition, e.g.  $R_j^i \lambda_i^k, R_i^k \lambda_i^k, \lambda_i^k \mu_k^j$ . Ordinary multiplication means numerical multiplication of the components indicated, and results therefore in another indexed symbol. The multiplications above would result respectively in symbols with indices as follows:

$$C_{jk}^{i\kappa}, \quad P_i^{k\kappa}, \quad \sigma_{ik\kappa}^{kj}$$

$$\text{where } C_{23}^{12} = R_2^1 \lambda_3^2, \quad C_{32}^{21} = R_3^2 \lambda_2^1$$

$$P_2^1 = R_2^1 \lambda_2^1, \quad P_3^3 = R_3^3 \lambda_3^3$$

$$\sigma_{13}^{12} = \lambda_4^1 \mu_3^{12}$$

etc. There are always as many variable indices in a product as there are different variable indices in the factors.

Thus  $R_j^i \lambda_i^k$  has three independent variable indices i, j, k and hence the product  $C_{jk}^{i\kappa}$  has three indices.

An important operation in indexed symbols is summation on one or more indices. This is so common in our work that

we indicate it by replacing the index in question by a large dot. Thus suppose the index h has a range of variation of 1 to 3 and i, a range 1 to 5, then

$$\lambda_{jk}^{hi} = \sum_{k=1}^3 \lambda_{jk}^{ki} = \lambda_{jk}^{1i} + \lambda_{jk}^{2i} + \lambda_{jk}^{3i}$$

$$\begin{aligned}\lambda_{jk}^{..} = \sum_{k=1}^3 \sum_{i=1}^5 \lambda_{jk}^{ki} &= \lambda_{jk}^{11} + \lambda_{jk}^{12} + \lambda_{jk}^{13} + \lambda_{jk}^{14} + \lambda_{jk}^{15} \\ &\quad + \lambda_{jk}^{21} + \lambda_{jk}^{22} + \lambda_{jk}^{23} + \lambda_{jk}^{24} + \lambda_{jk}^{25} \\ &\quad + \lambda_{jk}^{31} + \lambda_{jk}^{32} + \lambda_{jk}^{33} + \lambda_{jk}^{34} + \lambda_{jk}^{35}\end{aligned}$$

Most of our indexed symbols will represent populations. Suppose we are considering two different Mendelian factors. Let the first have two alleles,  $A_1$  and  $A_2$ , and suppose the second factor has three;  $B_1$ ,  $B_2$ , and  $B_3$ . Then any population may be divided into 21 genetically different groups, having the genetic formulae

- |                       |                        |                        |
|-----------------------|------------------------|------------------------|
| (1) $A_1 A_1 B_1 B_1$ | (7) $A_2 A_2 B_1 B_1$  | (13) $A_1 A_2 B_1 B_1$ |
| (2) $A_1 A_1 B_1 B_2$ | (8) $A_2 A_2 B_1 B_2$  | (14) $A_1 A_2 B_1 B_2$ |
| (3) $A_1 A_1 B_1 B_3$ | (9) $A_2 A_2 B_1 B_3$  | (15) $A_1 A_2 B_1 B_3$ |
| (4) $A_1 A_1 B_2 B_2$ | (10) $A_2 A_2 B_2 B_2$ | (16) $A_1 A_2 B_2 B_2$ |
| (5) $A_1 A_1 B_2 B_3$ | (11) $A_2 A_2 B_2 B_3$ | (17) $A_1 A_2 B_2 B_3$ |
| (6) $A_1 A_1 B_3 B_3$ | (12) $A_2 A_2 B_3 B_3$ | (18) $A_1 A_2 B_3 B_3$ |
|                       |                        | (19) $A_1 A_2 B_2 B_1$ |
|                       |                        | (20) $A_1 A_2 B_3 B_1$ |
|                       |                        | (21) $A_1 A_2 B_3 B_2$ |

This population would be represented by the symbol  $\lambda_{jk}^{hi}$ .

The indices  $h$  and  $j$  correspond to the first locus and since there are two alleles for this factor they each have a range of variation of 1 to 2. The second factor has three alleles and correspondingly  $i$  and  $k$  range over the values 1, 2, 3. Now, thinking of  $h$ ,  $i$ ,  $j$ , and  $k$  as fixed or dead we define the components of  $\lambda_{j,k}^{h,i}$  in the following manner. If  $h = j$  and  $i = k$  then  $\lambda_{j,k}^{h,i} = \lambda_{h,i}^{h,i} =$  the fraction of the population with the genetic formula  $A_h A_h B_i B_i$ . Thus  $\lambda_{1,3}^{1,3}$  is the fraction of the population of the type  $A_1 A_1 B_3 B_3$ . If  $h \neq j$ , or  $i \neq k$ , or both, then  $\lambda_{j,k}^{h,i}$  represents one half the fraction of the population having the formula  $A_h A_j B_i B_k$  or what is the same thing  $A_j A_h B_k B_i$ . Thus  $\lambda_{1,3}^{1,3}$  and  $\lambda_{1,2}^{1,3}$  are one half the fractions having the respective formulas  $A_1 A_2 B_3 B_3$  and  $A_2 A_1 B_3 B_2$ .

We shall use Greek letters as base letters for populations, and in general, then, the symbol

$$\lambda_{j_1, j_2, \dots, j_s}^{i_1, i_2, \dots, i_s} \quad (4)$$

represents a population in which s gene loci are under consideration. The first column of indices,  $i_1$  and  $j_1$ , corresponds to the first factor under consideration, the second column to the next factor, etc. Each factor may have an arbitrary number of alleles, and the linkage between any two may be of any value including 50% or random assortment. In case  $i_1=j_1, i_2=j_2, \dots, i_s=j_s$  then

$\lambda^{i_1 i_2 \dots i_s}_{i_1 i_2 \dots i_s}$  is the fraction of the population having the formula  $(A_{i_1} A_{i_2} B_{i_3} B_{i_4} \dots S_{i_s} S_{i_s})$ . If these equalities are not all true then

$\lambda^{i_1 i_2 \dots i_s}_{j_1 j_2 \dots j_s}$  is 1/2 the fraction having the formulae  $(A_{i_1} A_{j_2} B_{i_2} B_{j_3} \dots S_{i_s} S_{j_s})$ . It is helpful in using this notation to note the close connection between the two rows of letters in the indices and the two rows of genes in the chromosomes. The analogue is more than superficial, for we will later show how crossing over, say between the second and third loci, is connected, in this notation, with the symbol  $\lambda^{i_1 \cdot}_{\cdot \cdot j_1 \cdot}$ .

### III. Fundamental Theorems

There are two fundamental manipulation laws which we present as theorems because of their importance, although both are almost obvious.

#### Theorem I.

$$\lambda_{i_1 i_2 \dots i_s}^{j_1 j_2 \dots j_s} = \lambda_{j_1 j_2 \dots j_s}^{i_1 i_2 \dots i_s} \quad (5)$$

That is, inverting the upper and lower rows of indices of a population gives an identical population. This is evident from the meaning of the symbols, since a genetic formula  $(A_{i_1} A_{i_2} B_{j_1} B_{j_2} \dots S_{k_1} S_{k_2})$  is identical with the formula  $(A_{j_1} A_{j_2} B_{i_1} B_{i_2} \dots S_{k_1} S_{k_2})$ . This inversion of indices may be carried out independent of the location of the gene loci in question. However, if it is known that certain of the loci are in one chromosome pair, and none of the others are in this pair, further identities will hold. Namely, we may invert the indices corresponding to this chromosome pair and leave the others intact without changing the meaning of the symbol. Thus, if in the population  $\lambda_{k l m}^{h i j}$ , the first two loci are in one chromosome pair, and the third in another, we have:

$$\lambda_{k l m}^{h i j} = \lambda_{k l j}^{h i m} = \lambda_{h i j}^{k l m} = \lambda_{h i m}^{k l j}.$$

Theorem II.

$$\lambda_{\dots\dots\dots}^{h\dots i\dots k\dots} = 1 \quad (6)$$

That is summation of a population on all indices gives the result one. Obviously the sum of all the fractional parts of the population is unity. Now those parts which we have divided by 2 come into the summation (6) twice, corresponding to an inversion of the upper and lower indices. Thus with the population  $\lambda_{j,k}^{hi}$  the term  $\lambda_{i,j}^{hi}$  appears twice in the summation  $\lambda_{\dots\dots\dots}^{h\dots i\dots k\dots}$ , once as  $\lambda_{i,j}^{hi}$  and once as  $\lambda_{j,i}^{hi}$ , which are equal by Theorem I.

The significance of summation on a fewer number of indices is also of considerable importance. Consider the population  $\lambda_{j,k}^{hi}$ . Summation on  $k$  gives  $\lambda_{j,\cdot}^{hi}$ , a symbol with three variable or "live" indices. The reader may verify that if  $h \neq j$ ,  $\lambda_{j,\cdot}^{hi}$  represents one half the fraction of  $\lambda_{j,k}^{hi}$  which have the formula  $A_h A_j B_i -$  where the  $-$  may be any gene. If  $h = j$  then  $\lambda_{j,\cdot}^{hi}$  represents  $1/2$  the fraction having the formula  $A_h A_j B_i -$  plus one half the fraction of the type  $A_h A_h B_i B_i$ .

Summing on both  $k$  and  $i$  we have  $\lambda_{j,\cdot}^{hi}$  and this may be shown to have exactly the same meaning as  $\lambda_j^{hi}$  where  $h$  and  $j$  refer to the same gene locus in each case. That is, summing on a pair of vertical indices is equivalent to eliminating this locus from consideration. Summing on two horizontal indices,  $j$  and  $k$ , gives  $\lambda_{\cdot,\cdot}^{hi}$ , a two index

symbol, whose components are the fractions of all chromosomes in the population in which genes  $A_h$  and  $B_i$  both appear. Likewise  $\lambda_{\cdot \cdot}^{h \cdot}$  represents the fraction of all chromosome pairs in which  $A_h$  is in one and  $B_k$  in the other.

Summation on three indices  $\lambda_{\cdot \cdot \cdot}^{h \cdot \cdot}$  gives a one index symbol and its components are the gene frequencies of  $A_1, A_2, \dots$  in the general population. That is  $\lambda_{\cdot \cdot \cdot}^{h \cdot \cdot}$  is the sum of the fraction of homozygotes in  $A_1$  and half the heterozygotes having one  $A_1$  gene.

The reader will easily generalize these statements for more complex cases. For easy reference we summarize the above remarks in the following proposition.

Theorem III.

1. If  $i_1 = j_1, i_2 = j_2, \dots, i_s = j_s$ , then  $\lambda_{j_1 j_2 \dots j_s}^{i_1 i_2 \dots i_s}$  is one half the fraction of the population of the type  $A_{i_1} A_{i_2} \dots S_{i_s}$  plus half the fraction of the type  $A_{i_1} A_{i_2} \dots S_{i_s} S_{j_{s+1}}$ . If the conditions above do not hold then it represents one half the population of the type  $A_{i_1} A_{i_2} B_{j_3} B_{j_4} \dots S_{i_s}$ .

2. Summation on two vertical indices is equivalent to eliminating the corresponding locus from consideration.

3. Summing on one index in each column gives the fractions of the chromosome pairs in which alleles with the remaining indices appear.

4. The fraction of  $\lambda_{j_1 \dots j_s}^{i_1 \dots i_s}$  having at least one  $A_h$  gene is given by

$$2 \lambda_{\cdot \cdot \cdot \cdot \cdot}^{h \cdot \cdot \cdot \cdot \cdot} - \lambda_{\cdot \cdot \cdot \cdot \cdot}^{h \cdot \cdot \cdot \cdot \cdot}$$

In order to present a rigorous mathematical development it is convenient to consider symbols whose components are not all positive real numbers lying between zero and one. Of course, such a symbol cannot represent an actual group of individuals, but in some cases it is possible to solve problems using these symbols and get an actual population for the final answer. We shall speak of the symbol as a population in either case, but if the symbol does correspond to a possible group of individuals we shall say the population is realizable. If the components are all real numbers we will say the population is real. The use of unrealizable populations also adds elegance and generality to some of our later theorems.

We now introduce an operation between two populations to be known as cross multiplication and written:

$$\lambda_{j^k}^{hi} \times \mu_{j^k}^{hi}$$

This will be defined in such a way that the cross product of two realizable populations represents the expected offspring if the population in question are cross mated at random. We will first give the definition for the case of two linked factors. Let the probability of zero or an even number of crossovers between the two factors be  $p_0 = 1 - p_1$ , with  $p_1$ , then, the probability of an odd number (i.e the recombination value). We define the

cross product of  $\lambda_{j,k}^{h,i}, \mu_{j,k}^{h,i}$  as follows

$$\begin{aligned}\nu_{j,k}^{h,i} &= \lambda_{j,k}^{h,i} \times \mu_{j,k}^{h,i} \\ &= \frac{1}{2} [p_0 \lambda_{..}^{h,i} + p_1 \lambda_{..}^{h,i}] [p_0 \mu_{..}^{h,i} + p_1 \mu_{..}^{h,i}] \\ &\quad + \frac{1}{2} [p_0 \lambda_{..}^{i,k} + p_1 \lambda_{..}^{i,k}] [p_0 \mu_{..}^{i,k} + p_1 \mu_{..}^{i,k}]\end{aligned}\quad (8)$$

To prove that this represents the expected offspring population we use the idea of gene-pair frequencies, a generalization of the idea of gene frequencies. We wish to determine the probability that an offspring of mating  $\lambda$  with  $\mu$  will have the genetic formula  $A_h A_j B_i B_k$ . One way an offspring may obtain this formula is to get  $A_h B_i$  from  $\lambda$  and  $A_j B_k$  from  $\mu$ . This corresponds to the term

$$[p_0 \lambda_{..}^{h,i} + p_1 \lambda_{..}^{h,i}] [p_0 \mu_{..}^{i,k} + p_1 \mu_{..}^{i,k}] \quad (9)$$

in the equation (8). Each of the terms of (9) is a gene-pair frequency. The first is the frequency of the gene pair  $A_h B_i$  in the population  $\lambda$ . The second is the frequency of the gene-pair  $A_j B_k$  in  $\mu$ . The product is then the probability of obtaining an offspring  $A_h A_j B_i B_k$  by the method described.

Gene-pair frequency here means the frequency with which this pair of genes occurs together in the same chromosome after crossovers have taken place. The term

$$[p_0 \lambda_{..}^{h,i} + p_1 \lambda_{..}^{h,i}] \quad (10)$$

actually represents this because, from Theorem III,  $\lambda^{h,i}$  represents the frequency with which  $A_h$  and  $B_i$  appear together in the same chromosome before crossover, and  $p_0$  is the probability that there is no crossover between these factors (or at most an even number) so that  $p, \lambda^{h,i}$  is the probability of getting  $A_h$  and  $B_i$  together after crossovers if they start together. Likewise  $p, \lambda^{h,i}$  is the probability of  $A_h$  and  $B_i$  ending together when they start in opposite chromosomes of a pair. Since these two possibilities are mutually exclusive, and collectively exhaustive, their sum (10) represents the gene-pair frequency of  $A_h B_i$  in  $\lambda$  after crossovers. Similarly the second term of (9) is the gene-pair frequency of  $A_j B_k$  in  $\mu$  after crossover and the product of these two is the probability of an offspring getting the pair  $A_h B_i$  from  $\lambda$  and  $A_j B_k$  from  $\mu$ . The only other way for an offspring to get the formula  $A_h A_j B_i B_k$  is to get  $A_j B_k$  from  $\lambda$  and  $A_h B_i$  from  $\mu$ . This corresponds in exactly the same way to the second term of (8), and we may add the probabilities since the events are mutually exclusive. All that remains to be explained in (8) is the factor  $1/2$ . In case the equations  $h = j$ ,  $i = k$  do not both hold then for the components of the offspring population we want half of the fraction of individuals of this type in order to fit our previous definition of a population symbol, and hence this factor. If both these equalities are true then both terms

of (8) are identical and we may add getting

$$\omega_{h,i}^{h,i} = [p_0 \lambda_{..}^{h,i} + p_1 \lambda_{.i}^{h,i}] [p_0 \mu_{..}^{h,i} + p_1 \mu_{.i}^{h,i}]$$

which is what we get by the derivation above in this case, since the two "different" possibilities become identical under this restriction and therefore (8) holds for all values of the indices.

For three factors the defining equation of the cross product is

$$\begin{aligned}\omega_{k,l,m}^{h,i,j} &= \lambda_{k,l,m}^{h,i,j} \times \mu_{k,l,m}^{h,i,j} \\&= \frac{1}{2} [p_{00} \lambda_{...}^{h,i,j} + p_{01} \lambda_{..j}^{h,i,j} + p_{10} \lambda_{.ij}^{h,i,j} + p_{11} \lambda_{.i.j}^{h,i,j}] \\&\quad \cdot [p_{00} \mu_{...}^{h,i,j} + p_{01} \mu_{..j}^{h,i,j} + p_{10} \mu_{.ij}^{h,i,j} + p_{11} \mu_{.i.j}^{h,i,j}] \quad (11) \\&+ \frac{1}{2} [p_{00} \lambda_{...}^{h,i,j} + p_{01} \lambda_{..m}^{h,i,j} + p_{10} \lambda_{.im}^{h,i,j} + p_{11} \lambda_{.i.m}^{h,i,j}] \\&\quad \cdot [p_{00} \mu_{...}^{h,i,j} + p_{01} \mu_{..m}^{h,i,j} + p_{10} \mu_{.im}^{h,i,j} + p_{11} \mu_{.i.m}^{h,i,j}]\end{aligned}$$

In this equation  $p_{00}$  is the probability of an even number of crossovers between the first two genes and an even number between the second and third. If we wish to consider interference effects we cannot merely write  $p_{00} = p_0 q_0$  with  $p_0$  the probability of an even number of crosses between the first two loci and  $q_0$  that for the second and third, since the events are not independent. However defining  $q_0$  as the probability of an even number of crosses

between the second and third factors after it is known that an even number occurred between the first two, would make this valid. Similarly  $p_{01}$  is the probability that an even number of crosses occur between the first two factors and an odd number between the second two, etc. The method of formation of the formula is fairly obvious; note first that all permutations of 0, 1 are used on the coefficient  $p$ . Also a 1 corresponds to changing from one row of indices to another, while 0 corresponds to staying in the same row.

The proof of formula (11), and indeed the general case, is an easy generalization of the method used for two factors. It merely amounts to showing that a term such as

$$[p_{00} \lambda^{hij} + p_{01} \lambda^{hij} + p_{10} \lambda^{hij} + p_{11} \lambda^{hij}]$$

is the gene-triplet frequency for the set  $A_h B_i C_j$  after crossovers in population  $\lambda_{hij}^{hij}$ . This the reader will readily verify.

For  $n$  linked genes the expression for the cross product of two populations will take the form:

$$\omega_{j_1 j_2 \dots j_s}^{i_1 i_2 \dots i_s} = \lambda_{j_1 j_2 \dots j_s}^{i_1 i_2 \dots i_s} \times \mu_{j_1 j_2 \dots j_s}^{i_1 i_2 \dots i_s}$$

$$= \frac{1}{2} \left\{ [ p_{00} \dots \lambda^{i_1 i_2 \dots i_s} + p_{10} \dots \lambda^{i_1 i_2 \dots i_s} \right. \\ \left. + \dots + p_{11} \dots \lambda^{i_1 i_2 \dots i_s} ] \right.$$

• [same expression with  $\lambda$  replaced by  $\mu$  and  $i_1, i_2, \dots, i_s$  by  $j, j_2, \dots, j_s] \} + \frac{1}{2} \{ \text{ same pair of expressions with } \lambda \text{ and } \mu \text{ interchanged} \}.$

Although we have spoken throughout as though the factors under consideration were linked and in the same chromosome, this is not necessary. Suppose that in equation (11) the first two genes have a recombination value  $p_1 = 1 - p_0$ , and that they are located in a different chromosome from the third. Under these conditions it is easy to see that

$$p_{00} = p_{01} = 1/2p_0$$

$$p_{10} = p_{11} = 1/2p_1$$

Also we have

$$\lambda_{k \ell m}^{hij} = \lambda_{k \ell j}^{him}; \quad \mu_{k \ell m}^{hij} = \mu_{k \ell j}^{him}$$

From these equations (11) may be reduced to

$$\frac{1}{2} [ p_0 \lambda^{hij} + p_1 \lambda^{hij} ] [ p_0 \mu^{k \ell m} + p_1 \mu^{k \ell m} ] \\ + \frac{1}{2} [ p_0 \lambda^{k \ell m} + p_1 \lambda^{k \ell m} ] [ p_0 \mu^{hij} + p_1 \mu^{hij} ]$$

So that the independence of factors merely simplifies the situation.

In case all three factors are independent

$$p_{00} = p_{01} = p_{10} = p_{11} = 1/4$$

$$\text{and } \lambda_{iklm}^{lij} = \lambda_{kilm}^{lim} = \lambda_{klm}^{ilm} \text{ etc.}$$

so that (11) reduces to

$$\frac{1}{2} \lambda_{...}^{lij} u_{...}^{kilm} + \frac{1}{2} \lambda_{...}^{ilm} u_{...}^{lij}$$

We have proved, then, the fundamental

Theorem IV. The cross product of two realizable populations represents the expected offspring of these populations when cross mated at random.

Much of our work will now be the investigation of special cases of the general formulae given above. We note at once both from the mathematical definition and obvious genetic considerations the following

Theorem V.

$$\lambda_{j_1 \dots j_s}^{i_1 \dots i_s} \times \mu_{j_1 \dots j_s}^{i_1 \dots i_s} = \mu_{j_1 \dots j_s}^{i_1 \dots i_s} \times \lambda_{j_1 \dots j_s}^{i_1 \dots i_s}$$

That is to say, cross multiplication is commutative.

Let us now consider the case of a single factor, but with, however, any number of alleles. Then the cross product reduces to

$$v_j^i = \lambda_j^i \times \mu_j^i = \frac{1}{2} [\lambda_j^i \mu_j^i + \lambda_j^{\bar{i}} \mu_j^{\bar{i}}] \quad (14)$$

from which we get the proposition:

Theorem VI. If  $\lambda^i = \mu^i$  then  $\lambda_j^i \sigma_j^i = \mu_j^i \sigma_j^i$  (15)

In other words if two populations have the same gene frequencies for all alleles of a factor, then they will have the same breeding characteristics when cross-mated with another population.

Theorem VII. If  $v_j^i = \lambda_j^i \times \mu_j^i$  then  $v^i = \frac{1}{2}(\lambda^i + \mu^i)$  (16)

This follows immediately from (14) on summing through on the index  $j$  and noting by Theorem II that  $\lambda^i = \mu^i = 1$ .

This theorem shows that a gene frequency of a cross product is the arithmetic mean of the corresponding gene frequencies of the factors in the product.

We have already indicated in Part II how indexed symbols are added. When our symbols represent populations we shall consider addition only when the sum of the coefficients is unity. The purpose of this restriction is to keep all terms on an actual fractional basis and thus preserve the validity of our theorems. In general this causes little or no inconvenience, as merely dividing by the sum of the coefficients will always normalize in this sense.

We write

$$\lambda_{jk}^{hi} = R_1 \mu_{jk}^{hi} + R_2 v_{jk}^{hi} + \dots + R_n \sigma_{jk}^{hi} \quad (17)$$

where  $\sum_{i=1}^n R_i = 1$ , for the "sum" of the populations

$\mu, \nu, \dots, \sigma$  in the fractional proportions  $R_1, R_2, \dots, R_n$

Note that all terms of a sum must have the same indices (although sometimes an index may be changed in position in the same vertical column). This is part of a useful idea in indexed symbols known as index balance.

Index balance serves as a simple partial check on equations. If the indices do not balance in an equation, the equation is certainly wrong (in fact it is meaningless). The rules governing index balance for our work may be formulated as follows:

1. Each term in a sum must have the same indices.
2. Each side of an equation must have the same indices.
3. A product (ordinary or cross) has indices corresponding to each different live index appearing on any of the factors of the product. (See Part II).

Index balance applies only to live indices. There is no balance, for example, on the dead indices 0 and 1 on  $p_{01}$  in equation (11).

Addition of populations (17) is interpreted very simply as the population obtained by combining random samples of  $\mu, \nu, \dots, \sigma$  in the fractional proportions  $R_1, R_2, \dots, R_n$ .

Theorem VIII. Cross multiplication is distributive on addition, e.g.

$$\lambda_{j,k}^{k-i}(R_1\mu_{j,k}^{k-i} + R_2\nu_{j,k}^{k-i}) = R_1\lambda_{j,k}^{k-i}\mu_{j,k}^{k-i} + R_2\lambda_{j,k}^{k-i}\nu_{j,k}^{k-i} \quad (18)$$

We shall prove the theorem only for this simple case. The method of proof, however, is perfectly general and will apply with any number of indices and any number of terms in the sum. The left side of the equation is, by the definition (8):

$$\begin{aligned}
 & \frac{1}{2} [p_0 \lambda^{h,i} + p_1 \lambda^{h,i}] [p_0 (R_0 \mu^{j,k} + R_1 \nu^{j,k}) + p_1 (R_0 \mu^{j,k} + R_1 \nu^{j,k})] \\
 & + \frac{1}{2} [p_0 \lambda^{j,k} + p_1 \lambda^{j,k}] [p_0 (R_0 \mu^{h,i} + R_1 \nu^{h,i}) + p_1 (R_0 \mu^{h,i} + R_1 \nu^{h,i})] \\
 = & R_1 \left\{ \frac{1}{2} [p_0 \lambda^{h,i} + p_1 \lambda^{h,i}] [\mu^{j,k} + \nu^{j,k}] \right. \\
 & \left. + \frac{1}{2} [p_0 \lambda^{j,k} + p_1 \lambda^{j,k}] [\mu^{h,i} + \nu^{h,i}] \right\} \\
 & + R_2 \left\{ \frac{1}{2} [p_0 \lambda^{h,i} + p_1 \lambda^{h,i}] [\nu^{j,k} + \mu^{j,k}] \right. \\
 & \left. + \frac{1}{2} [p_0 \lambda^{j,k} + p_1 \lambda^{j,k}] [\nu^{h,i} + \mu^{h,i}] \right\} \\
 = & R_1 \lambda_{j,k}^{h,i} \times \mu_{j,k}^{h,i} + R_2 \lambda_{j,k}^{h,i} \times \nu_{j,k}^{h,i}.
 \end{aligned}$$

The theorem on equilibrium of population after random intermatting may be easily proved by the methods we have developed. We shall consider a somewhat more general case than is usually used, in that we allow any number of alleles, and also cross breeding between generations (i.e. the generations need not come in distinct steps nor need each individual mate in its own generation). We prove, then,

Theorem IX.

$$R_1 \lambda_i^k \cdot \lambda_i^k + R_2 \lambda_i^k \times (\lambda_i^k \times \lambda_i^k) + R_3 (\lambda_i^k \times \lambda_i^k) \times (\lambda_i^k \times \lambda_i^k)$$
$$+ \dots = \lambda_i^k \times \lambda_i^k = \lambda_i^k \lambda_i^k \quad (19)$$

The first term corresponds to a component representing direct offspring of our present population  $\lambda_i^k$  the second term represents a fraction obtained by mating this offspring with the parent generation, etc. Consider any term of this expression. In order to have a meaning it must have a factor  $(\lambda_i^k \times \lambda_i^k) = \lambda_i^k \lambda_i^k$  (by 14), but this may be replaced by  $\lambda_i^k$  from Theorem VI since  $\lambda_i^k \lambda_i^k = \lambda_i^k$ . Hence the number of factors in the product may be reduced by one. Continuing in this manner all terms reduce to the form  $R_n \lambda_i^k \times \lambda_i^k$  and adding we get the desired result.

In particular, if we have "step type" generations this result shows that

$$\begin{aligned} \lambda_i^k \times \lambda_i^k &= (\lambda_i^k \times \lambda_i^k) \times (\lambda_i^k \times \lambda_i^k) \\ &= [(\lambda_i^k \times \lambda_i^k) \times (\lambda_i^k \times \lambda_i^k)] \times [(\lambda_i^k \times \lambda_i^k) \times (\lambda_i^k \times \lambda_i^k)] \\ &= \text{etc.} \end{aligned}$$

these expressions being the 2nd, 3rd, 4th, etc. offspring generations, and all of these are equal to  $\lambda_i^k \lambda_i^k$ . It is obvious that a necessary and sufficient

condition for a population  $\lambda_i^h$  to be in this type of equilibrium is that

$$\lambda_i^h = \lambda^h \cdot \lambda^i$$

For two alleles it is well known that this is equivalent to the condition

$$\lambda^1 \lambda^2 = \lambda^1 \lambda^1$$

We now show how this, and the generalized result for any number of alleles may be obtained.

Theorem X. The three following statements are all equivalent.

(1)  $\lambda_i^h = \lambda^h \cdot \lambda^i = \lambda_i^h \cdot \lambda_i^h$

(2)  $\lambda^h \lambda^i = \lambda_i^h \lambda_i^h$

(3) The matrix  $\|\lambda_i^h\|$  is of rank one.

By the matrix  $\|\lambda_i^h\|$  is meant the matrix:

$$\begin{matrix} \lambda^1 & \lambda^2 & \dots & \lambda^s \\ \lambda^1 & \lambda^2 & \dots & \lambda^s \\ \vdots & \vdots & \ddots & \vdots \\ \lambda^1 & \lambda^2 & \dots & \lambda^s \end{matrix}$$

In the first place (1) implies (2), for if

$$\lambda_i^h = \lambda_+^h \lambda_-^i$$

then

$$\lambda_{+}^h = \lambda_+^h \lambda_-^h$$

and

$$\begin{aligned}\lambda_+^h \lambda_-^i &= \lambda_+^h \lambda_+^h \lambda_-^i \lambda_-^i \\ &= (\lambda_+^h \lambda_+^i)(\lambda_-^h \lambda_-^i) = \lambda_i^h \lambda_i^h\end{aligned}$$

Also (2) implies (1). Summing (2) on i gives:

$$\lambda_i^h = \lambda_+^h \lambda_-^i$$

Hence:

$$\begin{aligned}\lambda_+^h \lambda_-^i &= \sqrt{\lambda_+^h} \lambda_i^h \\ &= \sqrt{\lambda_i^h} \lambda_i^h = \lambda_i^h\end{aligned}$$

Thus (1) and (2) are equivalent.

Condition (3) is equivalent to either of these,  
for if (1) is true:

$$\lambda_i^h = \lambda_+^h \lambda_-^i$$

and the elements of the matrix  $\|\lambda_i^h\|$  can be written  
as the product of a number depending only on the row  
by a number depending only on the column. This is a  
well known condition that the matrix be of rank not  
greater than one. The rank is actually one since at  
least one element is different from zero to satisfy  
 $\lambda_+^h = 1$ . Thus (1) implies (3). If (3) is true then

each second order minor of  $\|\lambda_i^k\|$  must vanish. In particular we have

$$\begin{vmatrix} \lambda_1^k & \lambda_2^k \\ \lambda_2^k & \lambda_1^k \end{vmatrix} = 0$$

Hence

$$\lambda_1^k \lambda_2^k = \lambda_1^k \lambda_2^k$$

so that (3) implies (2). This shows that all the conditions are equivalent and proves the theorem.

If a population is in equilibrium we have

$$\sum_i \sqrt{\lambda_i^k} = \sum_i \sqrt{\lambda_i^k \lambda_i^k} = \sum_i \lambda_i^k = \lambda_i^k = 1$$

but this is not a sufficient condition for equilibrium, as the example

$$\|\lambda_i^k\| = \begin{vmatrix} \frac{1}{9} & \frac{1}{6} & 0 \\ \frac{1}{6} & \frac{1}{9} & \frac{1}{6} \\ 0 & \frac{1}{6} & \frac{1}{9} \end{vmatrix}$$

proves, for here  $\sqrt{\frac{1}{9}} + \sqrt{\frac{1}{9}} + \sqrt{\frac{1}{9}} = 1$  while  $\left| \begin{matrix} \frac{1}{6} & 0 \\ \frac{1}{9} & \frac{1}{6} \end{matrix} \right| = \frac{1}{36} \neq 0$

so the population is not in equilibrium.

In case more than one factor is considered the population will, in general, only reach equilibrium (for gene combinations) asymptotically. Suppose we have two linked factors and assume "step type" gene-

rations. The result of random intermating is given by:

Theorem XI. Under random intermating of  $\lambda_{j,k}^{l,i}$  the  $n$ th generation is the population

$$\begin{aligned} u_{j,k}^{l,i} &= [p_0^{n-1}(p_0 \lambda_{..}^{l,i} + p_1 \lambda_{..}^{l,i}) + (1-p_0^{n-1}) \lambda_{..}^l \lambda_{..}^i] \\ &\cdot [p_0^{n-1}(p_0 \lambda_{..}^{j,k} + p_1 \lambda_{..}^{j,k}) + (1-p_0^{n-1}) \lambda_{..}^j \lambda_{..}^k] \quad (20) \end{aligned}$$

and (assuming  $p_0 \neq 1$ ) approaches asymptotically the population

$$u_{j,k}^{l,i} = \lambda_{..}^l \lambda_{..}^{j,i} \lambda_{..}^j \lambda_{..}^k \quad \text{as } n \rightarrow \infty$$

Proof: By definition (8) the first generation is:

$$\begin{aligned} \sigma_{j,k}^{l,i} &= \lambda_{j,k}^{l,i} \times \lambda_{j,k}^{l,i} = [p_0 \lambda_{..}^{l,i} + p_1 \lambda_{..}^{l,i}] [p_0 \lambda_{..}^{j,k} + p_1 \lambda_{..}^{j,k}] \\ &= [(p_0 \lambda_{..}^{l,i} + p_1 \lambda_{..}^{l,i}) + (1-p_0^{l-1}) \lambda_{..}^l \lambda_{..}^{i,i}] \\ &\cdot [(p_0 \lambda_{..}^{j,k} + p_1 \lambda_{..}^{j,k}) + (1-p_0^{l-1}) \lambda_{..}^j \lambda_{..}^{i,k}] \end{aligned}$$

so the theorem is true for  $n=1$ . We now show that if it is true for the  $n$ th generation it will be true for the  $(n+1)$ th generation and thus complete the proof by mathematical induction. Assume, then, that the  $n$ th generation is

$$\begin{aligned} u_{j,k}^{l,i} &= [p_0^{n-1}(p_0 \lambda_{..}^{l,i} + p_1 \lambda_{..}^{l,i}) + (1-p_0^{n-1}) \lambda_{..}^l \lambda_{..}^{i,i}] \\ &\cdot [p_0^{n-1}(p_0 \lambda_{..}^{j,k} + p_1 \lambda_{..}^{j,k}) + (1-p_0^{n-1}) \lambda_{..}^j \lambda_{..}^{i,k}] \quad (21) \end{aligned}$$

whence, summing on  $j$  and  $k$ :

$$\begin{aligned} \mu_{..}^{k,i} &= [p_0^{n-1}(p_0 \lambda_{..}^{k,i} + p_1 \lambda_{..}^{k,i}) + (1-p_0^{n-1}) \lambda_{..}^{k,i} \lambda_{..}^{k,i}] \\ &\quad \cdot [p_0^{n-1}(p_0 \lambda_{..}^{k,i} + p_1 \lambda_{..}^{k,i}) + (1-p_0^{n-1}) \lambda_{..}^{k,i} \lambda_{..}^{k,i}] \\ &= [p_0^{n-1}(p_0 \lambda_{..}^{k,i} + p_1 \lambda_{..}^{k,i}) + (1-p_0^{n-1}) \lambda_{..}^{k,i} \lambda_{..}^{k,i}] \end{aligned}$$

since  $\lambda_{..}^{k,i} = 1$  and  $p_0 + p_1 = 1$ .

Also:

$$\begin{aligned} \mu_{..}^{k,k} &= [p_0^{n-1}(p_0 \lambda_{..}^{k,k} + p_1 \lambda_{..}^{k,k}) + (1-p_0^{n-1}) \lambda_{..}^{k,k} \lambda_{..}^{k,k}] \\ &\quad \cdot [p_0^{n-1}(p_0 \lambda_{..}^{k,k} + p_1 \lambda_{..}^{k,k}) + (1-p_0^{n-1}) \lambda_{..}^{k,k} \lambda_{..}^{k,k}] \\ &= \lambda_{..}^{k,k} \lambda_{..}^{k,k} \end{aligned}$$

Now the  $(n+1)$ th generation is given by the cross product of the  $n$ th generation with itself, i.e.

$$\begin{aligned} \mu_{..}^{k,k} \times \mu_{..}^{k,i} &= [p_0 \mu_{..}^{k,i} + p_1 \mu_{..}^{k,i}] [p_0 \mu_{..}^{k,k} + p_1 \mu_{..}^{k,k}] \\ &= [p_0 \{ p_0^{n-1}(p_0 \lambda_{..}^{k,i} + p_1 \lambda_{..}^{k,i}) + (1-p_0^{n-1}) \lambda_{..}^{k,i} \lambda_{..}^{k,i} \} + p_1 \lambda_{..}^{k,i} \lambda_{..}^{k,i}] \\ &\quad \cdot [p_0 \{ p_0^{n-1}(p_0 \lambda_{..}^{k,k} + p_1 \lambda_{..}^{k,k}) + (1-p_0^{n-1}) \lambda_{..}^{k,k} \lambda_{..}^{k,k} \} + p_1 \lambda_{..}^{k,k} \lambda_{..}^{k,k}] \\ &= [p_0^n(p_0 \lambda_{..}^{k,i} + p_1 \lambda_{..}^{k,i}) + (1-p_0^n) \lambda_{..}^{k,i} \lambda_{..}^{k,i}] \\ &\quad \cdot [p_0^n(p_0 \lambda_{..}^{k,k} + p_1 \lambda_{..}^{k,k}) + (1-p_0^n) \lambda_{..}^{k,k} \lambda_{..}^{k,k}] \end{aligned}$$

which is the same expression as (21) with  $n$  replaced by  $(n+1)$ . This, therefore, completes the proof. The asymptotic value is obvious since if  $p_0 \neq 1$  then as  $n \rightarrow \infty, p_0^{n-1} \rightarrow 0$  and  $(1-p_0^{n-1}) \rightarrow 1$  so that (21) reduces to

$$\lambda^h : \lambda^i : \lambda^j : \lambda^k \quad (22)$$

An obvious corollary is that a necessary and sufficient condition for a population  $C_{ijk}^{hi}$  to be stable under random intermatting is that it satisfy the conditions

$$C_{ijk}^{hi} = C^h : C^i : C^j : C^k$$

or that either  $p_0 = 0$  or  $p_0 = 1$ . If  $p_0 = 0$  the expression (21) reduces to its equilibrium value

$$\lambda^h : \lambda^i : \lambda^j : \lambda^k$$

at the first generation. If  $p_0 = 1$  we have perfect linkage and the expression becomes  $\lambda^h : \lambda^k$  as it should, since it then acts like a single factor.

We note that the speed with which equilibrium is approached depends entirely on the value of  $p_0$ . If  $p_0$  is small the approach will be very rapid, less so as  $p_0$  becomes larger, or the linkage closer between the factors.

It is interesting to see that with a given popula-

tion  $\lambda_{jk}^{ki}$ , the equilibrium will be approached more rapidly if there is very weak linkage ( $p_0 < 1/2$ ) than if the factors are completely independent, either in different chromosomes or in the same chromosome with  $p_0 = 1/2$ .

Incidentally, if there is no interference between crosses a recombination value  $p_0 < 1/2$  is mathematically impossible. Suppose the map distance (measured in morgans) between two loci is  $d$ . Let this distance be divided into a large number,  $n$ , of sections, each of map distance  $d/n$ . Then if  $n$  is large  $d/n$  is small and the probability of a cross in any one section is approximately  $d/n$  and approaches this value as  $n \rightarrow \infty$ . The probability of exactly  $s$  crosses between the two loci is given by

$$\lim_{n \rightarrow \infty} {}_n C_s \left(\frac{d}{n}\right)^s \left(1 - \frac{d}{n}\right)^{n-s} \quad (24)$$

Where  ${}_n C_s = \frac{n(n-1)\dots(n-s+1)}{s!}$  is the number of ways we can pick out the  $s$  sections where the crosses may occur,  $\left(\frac{d}{n}\right)^s$  is the probability that these crosses do occur and  $\left(1 - \frac{d}{n}\right)^{n-s}$  is the probability that crosses do not occur in the other  $(n-s)$  sections. This limit may be written as follows:

$$\lim_{n \rightarrow \infty} \frac{n(n-1)(n-2)\dots(n-s+1)}{n \cdot n \cdot n \dots n} \cdot \frac{d^s}{s!} \left(1 - \frac{d}{n}\right)^n \cdot \left(1 - \frac{d}{n}\right)^{-s}$$

The first factor approaches the limit 1, the second is a constant  $\frac{d^s}{s!}$ , the third is (setting  $x = \frac{n}{d}$ ):

$$\lim_{n \rightarrow \infty} \left(1 - \frac{d}{n}\right)^n = \lim_{x \rightarrow \infty} \left(1 - \frac{1}{x}\right)^{x^d} = e^{-d}$$

and the last term approaches 1 since the exponent is a constant. The entire function therefore approaches

$$\frac{d^s}{s!} e^{-d}$$

as the probability of exactly  $s$  crosses. The probability of an odd number is then

$$\begin{aligned} p_1 &= e^{-d} \left( \frac{d}{1!} + \frac{d^3}{3!} + \frac{d^5}{5!} + \dots \right) \\ &= e^{-d} \left( \frac{e^d - e^{-d}}{2} \right) = \frac{1 - e^{-2d}}{2} \end{aligned}$$

which is clearly less than 50% for any real  $d$ . In case interference is present we cannot multiply probabilities as we did in (24) since the events are no longer independent and it is at least mathematically possible for values of  $p_1 > 1/2$ . Thus suppose we have a long chromosome which is very likely to cross over at least once, but one cross strongly inhibits any other crosses for a large distance. It is evident that such conditions would allow recombination values greater than 50%.

In the case when the factors are in different chromosomes, equation (20) can be simplified, due to the fact that

$$\lambda_{j,k}^{i,i} = \lambda_{j,i}^{k,k}$$

under these conditions. It follows that

$$\lambda_{..}^{h,i} = \lambda_{..}^{h,i}; \quad \lambda_{..}^{j,k} = \lambda_{..}^{j,k}$$

so that (20) reduces to

$$\begin{aligned} \mu_{j,k}^{h,i} &= [p_0^{n-1} \lambda_{..}^{h,i} + (1-p_0^{n-1}) \lambda_{..}^{h..} \lambda_{..}^{i..}] \\ &\quad [p_0^{n-1} \lambda_{..}^{j,k} + (1-p_0^{n-1}) \lambda_{..}^{j..} \lambda_{..}^{k..}] \end{aligned}$$

For three linked factors the first generation offspring of the population  $\lambda_{k,l,m}^{h,i,j}$  is

$$\begin{aligned} &[p_{00} \lambda_{..}^{h,i,j} + p_{01} \lambda_{..}^{h,i..} + p_{10} \lambda_{..}^{h..i} + p_{11} \lambda_{..}^{h..j}] \\ &[p_{00} \lambda_{..}^{k,l,m} + p_{01} \lambda_{..}^{k,l..} + p_{10} \lambda_{..}^{k..l} + p_{11} \lambda_{..}^{k..m}] \quad (23) \end{aligned}$$

The second random intermating gives:

$$\begin{aligned} &[p_{00}(p_{00} \lambda_{..}^{h,i,j} + p_{01} \lambda_{..}^{h,i..} + p_{10} \lambda_{..}^{h..i} + p_{11} \lambda_{..}^{h..j}) \\ &+ p_{01}((p_{00} + p_{01}) \lambda_{..}^{h..i} + (p_{01} + p_{11}) \lambda_{..}^{h..j}) \lambda_{..}^{i..j} \\ &+ p_{10}((p_{00} + p_{10}) \lambda_{..}^{i..j} + (p_{01} + p_{11}) \lambda_{..}^{h..j}) \lambda_{..}^{h..i} \quad (25) \\ &+ p_{11}((p_{00} + p_{11}) \lambda_{..}^{h..j} + (p_{01} + p_{10}) \lambda_{..}^{h..i}) \lambda_{..}^{i..j}] \\ &\quad [\text{same expression with } h, i, j \text{ replaced by } k, l, m.] \end{aligned}$$

The nth generation has also been determined in

this case. It is given by the following proposition.

Theorem XII. Under random intermating of the population  $\lambda_{k l m}^{h i j}$  the  $n$ th offspring generation is given by:

$$\begin{aligned}\lambda_{k l m}^{h i j} = & \left[ p_{00}^{n-1} (p_{00} \lambda_{...}^{h i j} + p_{01} \lambda_{...}^{h i j} + p_{10} \lambda_{...}^{h i j} + p_{11} \lambda_{...}^{h i j}) \right. \\ & + ((p_{00} + p_{01})^{n-1} - p_{00}^{n-1}) \left\{ (p_{00} + p_{01}) \lambda_{...}^{h i j} + (p_{10} + p_{11}) \lambda_{...}^{h i j} \right\} \lambda_{...}^{h i j} \\ & + ((p_{00} + p_{10})^{n-1} - p_{00}^{n-1}) \left\{ (p_{00} + p_{10}) \lambda_{...}^{h i j} + (p_{10} + p_{11}) \lambda_{...}^{h i j} \right\} \lambda_{...}^{h i j} \\ & + ((p_{00} + p_{11})^{n-1} - p_{00}^{n-1}) \left\{ (p_{00} + p_{11}) \lambda_{...}^{h i j} + (p_{01} + p_{10}) \lambda_{...}^{h i j} \right\} \lambda_{...}^{h i j} \\ & \left. + (1 + 2p_{00}^{n-1} - (p_{00} + p_{01})^{n-1} - (p_{10} + p_{00})^{n-1} - (p_{00} + p_{11})^{n-1}) \lambda_{...}^{h i j} \lambda_{...}^{h i j} \lambda_{...}^{h i j} \right] \\ & \cdot \left[ \text{same expression with } h, i, j \text{ replaced by } k, l, m \right]\end{aligned}$$

This may be proved by mathematical induction exactly as we proved Theorem XI. The expression

approaches  $\lambda^h \dots \lambda^i \dots \lambda^j \dots \lambda^k \dots \lambda^l \dots \lambda^m$  asymptotically as  $n \rightarrow \infty$ . In case the three loci are in different chromosomes it reduces to

$$\frac{1}{16^{n-1}} \left[ \lambda^{h+i} + (2^{n-1}-1) \left\{ \lambda^{h+i} \lambda^{i+j} + \lambda^{h+j} \lambda^{i+h} \right. \right. \\ \left. \left. + \lambda^{i+j} \lambda^{h+i} \right\} + (4^{n-1}-3 \cdot 2^{n-1}+2) \lambda^{h+i} \lambda^{i+j} \lambda^{j+h} \right] \quad (26)$$

[same expression with h, i, j replaced by k, l, m]

It is possible to expand a population in a series form which displays the homogeneous components of the population. This series is very similar to the expansion of a Boolean function in Symbolic Logic, and not only throws light on the mathematical nature of the symbols we are using, but is also useful for computational purposes. To develop the expansion we must first define a set of "constants"; homogeneous populations of a certain fixed genetic constitution. These constants will always be represented by the base letter  $\gamma$ , and the indices refer to the particular locus or loci we are considering. All the members of a constant  $\gamma$  population have the same genetic constitution with respect to the factors under consideration.

For a single locus our definitions are as follows:

$$\alpha \gamma_i^j = \begin{cases} 1 & \text{if } i = j = a \\ 0 & \text{otherwise} \end{cases} \quad (27)$$

If  $a \neq b$

$$\alpha \gamma_i^j = \begin{cases} 1 & \text{if } i = a \text{ and } j = b \text{ or if } i = b \text{ and } j = a \\ 0 & \text{otherwise} \end{cases} \quad (28)$$

Here  $a$  and  $b$  are dead indices; they represent certain fixed numbers while the live indices  $i$  and  $j$  represent any of several values. The dead indices merely serve to distinguish one  $\gamma$  from another and there is, in general, no index balance on them. It will be seen from our definition that  $\gamma_j^i$  represents a population whose members are all of the genetic type  $A_a A_b$ . Thus  $\gamma_j^i$  is a homogeneous population with the formula  $A_1 A_2$ .

If we are considering only a single dimorphic factor the series expansion of a population  $\lambda_j^i$  is as follows:

$$\lambda_j^i = \lambda'_i \gamma_j^i + 2\lambda''_i \gamma_j^i + \lambda'''_i \gamma_j^i \quad (29)$$

To prove this it is merely necessary to note that it reduces to an identity for all values of  $i$  and  $j$ . Thus with  $i=1$ ,  $j=1$  all the terms on the right are zero, except the first which reduces to  $\lambda'_i$ . For  $i=1$ ,  $j=2$  only the second term is effective giving  $\lambda''_i = 2\lambda''_{i1} = \lambda''_i$ , etc.

The expansion (29) displays  $\lambda_j^i$  as a population made up of three homogeneous parts  $\gamma_j^i$ ,  $\gamma_j^i$ ,  $\gamma_j^i$ .

The fraction of each type is the coefficient of the corresponding  $\gamma$  in the expansion. For more than two alleles the expansion takes the form:

$$\begin{aligned}\lambda_j^i &= \lambda_1^i \gamma_j^i + \lambda_2^i \gamma_j^i + \dots + \lambda_n^i \gamma_j^i \\ &\quad + 2\lambda_1^i \gamma_j^i + 2\lambda_2^i \gamma_j^i + \dots + 2\lambda_n^i \gamma_j^i \\ &\quad + 2\lambda_3^i \gamma_j^i + \dots + 2\lambda_{n-1}^i \gamma_j^i\end{aligned}\tag{50}$$

With more than one factor under consideration we define the  $\gamma$  populations according to the following scheme:

$$\begin{aligned}\text{def } \gamma_{j,k}^{hi} &= \begin{cases} 1 & \text{if } \begin{cases} h = j = a \\ i = k = b \end{cases} \\ 0 & \text{otherwise} \end{cases} \\ &\quad \text{otherwise}\end{cases}\end{aligned}\tag{31}$$

If  $a \neq c$  or  $b \neq d$  or both

$$\begin{aligned}\text{def } \gamma_{j,k}^{hi} &= \begin{cases} 1/2 & \text{if } \begin{cases} h = a \quad i = b \\ j = c \quad k = d \end{cases} \text{ or } \begin{cases} h = c \quad i = d \\ j = a \quad k = b \end{cases} \\ 0 & \text{otherwise} \end{cases}\end{aligned}\tag{32}$$

Thus  $\gamma_{j,k}^{hi}$  represents a homogeneous population whose members all have the formula  $A_1 A_3 B_2 B_2$ . Constants for more than two factors are defined in a completely analogous manner. Thus  $\gamma_{j,k,m}^{hi,j}$  is a population whose members are of the type  $A_a A_d B_b B_c C_c C_f$ .

The series expansion for more than one locus has the form (taking three factors to be specific):

$$\lambda_{xlm}^{hij} = \lambda_{,,,}^{,,,} \gamma_{klm}^{hij} + 2\lambda_{,,}^{,,} \gamma_{klm}^{hij} \\ + \dots + \lambda_{r_1 r_2 r_3}^{r_1 r_2 r_3} \gamma_{r_1 r_2 r_3 klm}^{hij} \quad (53)$$

where  $r_1$ ,  $r_2$ , and  $r_3$  are the number of alleles of the three factors. Each term corresponding to a part of the population homozygous in all factors has the coefficient one; if the corresponding part of the population is heterozygous in one or more factors the coefficient should be two.

The cross product of any two  $\gamma$  population may be written as a linear combination of  $\gamma$ 's. Thus

$$c\gamma_i^h c\gamma_i^l = \frac{1}{4} c\gamma_i^h + \frac{1}{4} d\gamma_i^h + \frac{1}{4} c\gamma_i^l + \frac{1}{4} d\gamma_i^l \quad (34)$$

In case some of the numbers  $a$ ,  $b$ ,  $c$ ,  $d$  are equal this expression is still true but may be simplified. Thus

if  $a = b$

$$a\gamma_i^h \times d\gamma_i^l = \frac{1}{2} c\gamma_i^h + \frac{1}{2} d\gamma_i^l$$

For two loci the law of multiplication of the  $\gamma$ 's is

$$\begin{aligned} \frac{ab}{cd} \gamma_{j^k}^{hi} \times \frac{ef}{gh} \gamma_{j^k}^{hi} &= \frac{p_0}{4} [a\gamma_{j^k}^{hi} + b\gamma_{j^k}^{hi} + c\gamma_{j^k}^{hi} + d\gamma_{j^k}^{hi}] \\ &+ \frac{p_0 p_1}{4} [a\gamma_{j^k}^{hi} + b\gamma_{j^k}^{hi} + c\gamma_{j^k}^{hi} + d\gamma_{j^k}^{hi} + e\gamma_{j^k}^{hi} + f\gamma_{j^k}^{hi} + g\gamma_{j^k}^{hi} + h\gamma_{j^k}^{hi}] \\ &+ \frac{c\gamma_{j^k}^{hi}}{ef} + \frac{d\gamma_{j^k}^{hi}}{gh} + \frac{e\gamma_{j^k}^{hi}}{ch} + \frac{f\gamma_{j^k}^{hi}}{gh} + \frac{g\gamma_{j^k}^{hi}}{ch} + \frac{h\gamma_{j^k}^{hi}}{gh} \\ &+ \frac{e\gamma_{j^k}^{hi}}{ef} + \frac{f\gamma_{j^k}^{hi}}{gh} + \frac{g\gamma_{j^k}^{hi}}{ch} + \frac{h\gamma_{j^k}^{hi}}{gh} \end{aligned}$$

The series expansion and law of multiplication of the  $\gamma$ 's display a population as a hypercomplex number, i.e. as a symbol of the form  $(a_1 e_1 + a_2 e_2 + \dots + a_n e_n)$  where the coefficients  $a_1, a_2 \dots a_n$  are numbers and the symbols  $e_1, e_2 \dots e_n$  are "unit vectors" with some given law of multiplication such that the product of any two of the  $e$ 's may be written as a linear combination of  $e$ 's.

It is well known that except for trivial cases and the case of ordinary complex numbers, no law of multiplication preserves all the commutative, associative and distributive laws of ordinary numbers. In our case the associative law of multiplication is sacrificed. Thus the product  $(\gamma_i^h \times \gamma_i^h) \times \gamma_i^h$  is, by (34):

$$(\gamma_i^h \times \gamma_i^h) \times \gamma_i^h = \gamma_i^h \times \gamma_i^h = \frac{1}{2} \gamma_i^h + \frac{1}{2} \gamma_i^h$$

while on the other hand

$$\begin{aligned} \gamma_i^h \times (\gamma_i^h \times \gamma_i^h) &= \gamma_i^h \times \left( \frac{1}{2} \gamma_i^h + \frac{1}{2} \gamma_i^h \right) \\ &= \frac{1}{2} \gamma_i^h + \frac{1}{2} \left( \frac{1}{2} \gamma_i^h + \frac{1}{2} \gamma_i^h \right) \\ &= \frac{3}{4} \gamma_i^h + \frac{1}{4} \gamma_i^h \end{aligned}$$

This is a simple example of a multiplication in which the associative law does not hold, and shows that in a cross product of several factors it is essential that parentheses be retained to indicate the order in which

multiplication is performed.

The series expansion (33) of a population shows how an arbitrary population may be written as the sum of a set of particular homogeneous populations, the  $\gamma$ 's. The choice of this particular set of populations as components was a matter of convenience, not of necessity. We now show that any set of populations satisfying a certain simple condition would do as well.

Consider a set of  $n$  populations  $\lambda, \mu, \dots, \sigma$ .

We omit writing the indices, but there may be any number of loci. We will say that these populations are linearly independent if there is no set of numbers  $a_1, a_2, \dots, a_n$ , not all zero, and that

$$a_1\lambda + a_2\mu + \dots + a_n\sigma = 0 \quad (35)$$

for all values of the live indices.

Theorem XIII: Any population  $\phi$  may be expressed uniquely as a linear combination of  $n$  linearly independent populations where  $n$  is the number of different possible genetic formulae for the factors considered.

To prove this, note that a necessary and sufficient condition that (35) have no solution for the  $a$ 's (not all zero) is that the determinant

$$\begin{vmatrix} \lambda_{11\dots 1} & \mu_{11\dots 1} & \dots & \sigma_{11\dots 1} \\ \lambda_{11\dots 2} & \mu_{11\dots 2} & \dots & \sigma_{11\dots 2} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_{n_1 n_2 \dots n_s} & \mu_{n_1 n_2 \dots n_s} & \dots & \sigma_{n_1 n_2 \dots n_s} \end{vmatrix} \quad (36)$$

be different from zero. In this determinant each population takes on the values of all components in a column; i.e. the values obtained by giving the indices all possible values.

Now the non-vanishing of (36) is also a necessary and sufficient condition for the existence of a unique solution for the  $b$ 's in the equations:

$$\Phi = b_1 \lambda + b_2 \mu + \dots + b_n \sigma$$

and this proves the theorem.

In passing we note that if we have a linked factors with  $r_1, r_2, \dots, r_s$  alleles respectively, then  $n$ , the number of different components is given by

$$n = \frac{r_1 r_2 \dots r_s (r_1 r_2 \dots r_s + 1)}{2} \quad (37)$$

We may think of  $2s$  positions in which genes may be placed. There are  $r_1$ , possibilities for the first and second positions,  $r_2$  for the third and fourth, etc., and therefore a total of  $r_1^2 r_2^2 r_3^2 \dots r_s^2$ . However, as an interchange of

the two chromosomes does not affect the genetic constitution we should divide this by two, except for the ones which are homozygous in all factors and were not counted twice. There are  $r_1 r_2 \dots r_s$  types of fully homozygous individuals and we may correct our formula, then, by adding this and then dividing by two:

$$n = \frac{r_1^2 r_2^2 \dots r_s^2 + r_1 r_2 \dots r_s}{2} = 1/2 [r_1 r_2 \dots r_s (r_1 r_2 \dots r_s + 1)]$$

In case the loci are not all in the same chromosome but spread out in a number of different ones, we may evaluate the expression (37) for each chromosome involved and multiply these results.

#### IV. The Solution of Equations Involving Unknown Populations

It is easy to write down equations of various types involving unknown populations. Many of these may be solved for the unknowns in terms of the known populations by means of the theorems we have developed. In general an equation represents some breeding experiment involving a population of unknown genetic constitution resulting in a genetically known population. In the following we shall use the letters  $\varphi, \psi, \chi \dots$  for base letters in unknown populations and  $\lambda, \mu, \nu \dots$  for known populations.

The general method of attack on these problems may be outlined as follows:

1. By summing on various indices we are able to evaluate gene frequencies, gene pair frequencies, etc. for the unknown populations. This ordinarily involves no more than the solution of one or more linear algebraic equations.

2. Knowing these we can evaluate cross products in which the unknowns appear, since a cross product depends only on the values of the population symbol with half the indices dotted.

3. With only linear terms remaining it is usually easy to solve the equations by or-

dinary methods for algebraic equations.

To illustrate how this is done we shall consider several examples. Suppose first, for simplicity, that only one locus is involved, and that we have the equation:

$$R_1 \varphi_i^k + R_2 \varphi_i^k \times \varphi_i^k + R_3 \varphi_i^k \times (\varphi_i^k \times \varphi_i^k) + \dots = \lambda_i^k \quad \sum R_i = 1 \quad (38)$$

with  $\varphi_i^k$  unknown and  $\lambda_i^k$  known.

By Theorem IX this reduces immediately to the form

$$R \varphi_i^k + S \varphi_i^k \times \varphi_i^k = \lambda_i^k ; \quad R + S = 1 \quad (39)$$

Summing on  $i$  we have by Theorem VII

$$R \varphi^k + \frac{1}{2} S \varphi^k \times \varphi^k + \frac{1}{2} S \varphi^k = \lambda^k$$

or  $\varphi^k = \lambda^k$

Hence by Theorem VI we may replace  $\varphi_i^k$  by  $\lambda_i^k$  in any product. Returning then to equation (39) we have

$$R \varphi_i^k + S \lambda_i^k \times \lambda_i^k = \lambda_i^k$$

$$\varphi_i^k = \frac{1}{R} \lambda_i^k - \frac{S}{R} \lambda_i^k \times \lambda_i^k$$

and this must be the unique solution of the equation, if a solution exists. To prove that is a solution we merely try it in the equation and find that it is satisfied.

A more general equation in one unknown is the following:

$$R \varphi_i^k (\mu_i^k \times \varphi_i^k) + S \varphi_i^k \times \varphi_i^k + T \varphi_i^k = \lambda_i^k \quad (40)$$

Summing on i:

$$R \left( \frac{1}{2} \varphi^k + \frac{1}{4} \mu^k + \frac{1}{4} \varphi^k \right) + S \varphi^k + T \varphi^k = \lambda^k$$

$$\begin{aligned} \varphi^k &= \frac{4}{3R+4S+4T} \lambda^k - \frac{R}{3R+4S+4T} \mu^k \\ &= \frac{4}{4-R} \lambda^k - \frac{R}{4-R} \mu^k \end{aligned}$$

Replacing  $\varphi_i^k$  in each product of (39) by the expression

$$\frac{4}{4-R} \lambda^k - \frac{R}{4-R} \mu^k \text{ gives:}$$

$$\begin{aligned} \varphi_i^k &= \frac{1}{T} \lambda_i^k - \left( \frac{4}{4-R} \lambda_i^k - \frac{R}{4-R} \mu_i^k \right) \cdot \\ &\quad \times \left[ \frac{4}{4-R} \lambda_i^k \times \mu_i^k + \frac{R^2}{4-R} \mu_i^k \times \mu_i^k + \frac{4S}{4-R} \lambda_i^k - \frac{RS}{4-R} \mu_i^k \right] \end{aligned}$$

as the solution of (39).

It may be easily shown that the above method is applicable to any single equation in one unknown  $\varphi_i^k$  providing the coefficient of  $\varphi_i^k$  does not equal zero. Such an equation always has a unique solution, although this solution may not always represent a realizable population.

A system of linear simultaneous equations may be solved by the ordinary methods for algebraic equations, since by fixing the indices we actually have such a sys-

tem of linear algebraic equations. Thus suppose

$$\sum_{i=1}^{\infty} R_{ij} \cdot \varphi_i^k = j \lambda_e^k \quad j = 1, 2, \dots, n \quad (41)$$

represents the system of equation with  $\varphi_1^k, \varphi_2^k, \dots, \varphi_n^k$  unknowns. The indices i and j here serve to distinguish between the different  $\varphi$  and  $\lambda$  populations. The solution of this system is

$$i \varphi_i^k = \sum_{j=1}^n \frac{M_{ji}}{|R_{ji}|} j \lambda_e^k \quad (42)$$

where  $|R_{ji}|$ , the determinant of the coefficients  $\neq 0$   $\neq 0$ , and  $M_{ji}$  is the cofactor of  $R_{ji}$ .

The reader may verify that simultaneous systems of equations with cross products involving the unknown populations may also be solved by these methods. Turning now to problems involving more than one gene locus the situation becomes a bit more complicated. Suppose again we have the equation in one unknown

$$R \varphi_j^{ki} \cdot \varphi_{jk}^{ki} + s \varphi_{ik}^{ki} = \lambda_j^{ki} \quad (43)$$

Referring to the definition of a cross product for two factors we see that knowing the two quantities  $\varphi_j^{ki}$  and  $\varphi_{jk}^{ki}$  the cross product in (43) is completely determined. Summing, then, on j and k we get

$$R \varphi_{ik}^{ki} + s \varphi_{ik}^{ki} = \lambda_{ik}^{ki} \quad (44)$$

and we see that to determine  $\varphi_{ik}^{ki}$  we should find  $\varphi_{ik}^{ki}$ .

Summing (45) on  $i$  gives us

$$\varphi_{..}^{hi} = \lambda_{..}^{hi}$$

similarly

$$\varphi_{..}^{li} = \lambda_{..}^{li}$$

hence

$$R \lambda_{..}^{li} \lambda_{..}^{hi} + s \varphi_{..}^{li} = \lambda_{..}^{li}$$
$$\varphi_{..}^{li} = \frac{1}{s} \lambda_{..}^{li} - \frac{R}{s} \lambda_{..}^{li} \lambda_{..}^{hi} \quad (46)$$

Substituting in (44)

$$\varphi_{..}^{li} = \frac{1}{s+Rp} (\lambda_{..}^{li} - \frac{Rp}{s} \lambda_{..}^{li} + R \lambda_{..}^{li} \lambda_{..}^{hi}) \quad (47)$$

Now from our original equation (43):

$$\varphi_{jk}^{li} = \frac{1}{s} \lambda_{jk}^{li} - \frac{R}{s} \varphi_{jk}^{li} \times \varphi_{jk}^{li}$$

and the cross product on the right may be calculated from the equations (46) and (47). This, then, gives the unique solution to the problem.

Generalizing these methods to systems of simultaneous equations in more than one unknown and with any number of gene loci is not difficult. Of course, the equations become larger and more cumbersome, but no new theoretical difficulties appear.

## V. Lethal Factors and Selection

The results of a selective action may be calculated by the methods we have developed. Suppose the chances of survival of the types  $A_h A_i$  are  $R_i^h$  where  $h$  and  $i$  take values over all alleles. Then starting with a population  $\lambda_i^h$  the population reaching maturity will be

$$\frac{1}{D_1} R_i^h \lambda_i^h$$

where  $D_1$  is a normalizing factor given by

$$D_1 = \sum_k \sum_i R_i^k \lambda_i^k \quad (48)$$

Sums on more than one index are of frequent occurrence in selection work and we adopt the convention that summation on two or more indices simultaneously will be indicated by placing a bar over these indices. Thus (48) would be written

$$D_1 = R_{\bar{i}}^{\bar{k}} \lambda_{\bar{i}}^{\bar{k}}$$

The first generation offspring of  $\lambda_i^h$  would be

$$\begin{aligned} u_{\bar{i}}^{\bar{k}} &= \frac{1}{D_1} (R_i^h \lambda_i^h) \times (R_{\bar{i}}^{\bar{k}} \lambda_{\bar{i}}^{\bar{k}}) \\ &= \frac{1}{D_1} (R_{\bar{i}}^{\bar{k}} \lambda_{\bar{i}}^{\bar{k}} R_i^h \lambda_i^h) \end{aligned}$$

and if

$$D_2 = D_1^2 R_{\bar{i}}^{\bar{k}} u_{\bar{i}}^{\bar{k}}$$

then

$$\frac{1}{D_2} R_i^h R_{\bar{i}}^{\bar{k}} \lambda_i^h R_{\bar{i}}^{\bar{k}} \lambda_{\bar{i}}^{\bar{k}}$$

will reach maturity.

The next generation is given by

$$\frac{1}{D_2} R_j^k R_s^k \lambda_s^k R_j^{\bar{k}} \lambda_j^{\bar{k}} R_i^{\bar{k}} R_s^{\bar{k}} \lambda_s^{\bar{k}} R_i^{\bar{k}} \lambda_i^{\bar{k}}$$

etc.

A population will be in equilibrium if and only if the equation

$$\lambda_i^k = \frac{1}{D} R_i^k \lambda^k \lambda^{\bar{k}} \quad (49)$$

is satisfied. This requires that

$$\lambda^k = \frac{1}{D} R_i^k \lambda^k \lambda^{\bar{k}} \quad (50)$$

or

$$\lambda^k \left( \frac{1}{D} R_i^k \lambda^{\bar{k}} - 1 \right) = 0 \quad (51)$$

This may be satisfied in two different ways. If none

of the  $\lambda_i^k = 0$ , we must have

$$R_i^k \lambda^{\bar{k}} = D \quad (52)$$

a system of linear algebraic equations with the unique solution (providing  $|R_i^k| \neq 0$ )

$$\lambda^k = \kappa \begin{vmatrix} R_1^k R_2^k \dots R_{k-1}^k & | & R_{k+1}^k \dots R_n^k \\ R_1^{\bar{k}} R_2^{\bar{k}} \dots R_{k-1}^{\bar{k}} & | & R_{k+1}^{\bar{k}} \dots R_n^{\bar{k}} \\ \vdots & \ddots & \vdots \\ R_1^n R_2^n \dots R_{k-1}^n & | & R_{k+1}^n \dots R_n^n \end{vmatrix} \quad (53)$$

where  $k$  is a constant determined by the condition

$\lambda_i = 1$  Equation (51) may also be satisfied if some of the  $\lambda_i = 0$ . We may not divide though by these components, but the remainder gives us a set of, say,  $m < n$  equations in the  $m$  nonvanishing components which will, in general, have a unique solution.

A population satisfying (49) above would be in equilibrium in the sense that the first generation expected offspring would be of the same genetic constitution. However, this equilibrium may not be stable, for the actual offspring will in general deviate somewhat from the expected offspring and the population will be stable if and only if this deviation tends to cause the next generation to return to the equilibrium position. The situation may be likened to a ball which may be either balanced on the top of a hill or placed at the lowest point in a valley. In either case the ball is in equilibrium, but only in the valley is it stable, for if given a slight displacement on the hill the ball will tend to run down, while in the valley it tends to return to the lowest point.

Although a complete set of necessary and sufficient conditions for stability of a population have not been found, we have the following proposition:

Theorem XIV. A necessary condition for stability of a realizable equilibrium population  $\lambda_i^*$  (no gene

frequency equals zero) under the selective action  $R_i^k$  is that

$$R_i^k < R_s^k \quad s \neq k \quad (54)$$

Proof: Let the  $R_i^k$  coefficients be multiplied by such a constant that the equilibrium population satisfies the normalized equation:

$$\lambda_i^k = R_i^k \lambda_i^k \lambda_i^{\bar{k}} \quad (55)$$

or, since no  $\lambda_i^k = 0$

$$R_i^k \lambda_i^{\bar{k}} = 1 \quad (56)$$

Let this population take on a small increment  $\Delta \lambda_i^k$  with  $\Delta \lambda_i^{\bar{k}} = 0$ . The result of one generation random intermixing of this displaced population  $\lambda_i^k + \Delta \lambda_i^k$  is

$$u_i^k = \frac{1}{D} R_i^k (\lambda_i^k + \Delta \lambda_i^k)(\lambda_i^{\bar{k}} + \Delta \lambda_i^{\bar{k}})$$

Whence

$$\begin{aligned} u_i^k &= \frac{1}{D} R_i^k (\lambda_i^k + \Delta \lambda_i^k)(\lambda_i^{\bar{k}} + \Delta \lambda_i^{\bar{k}}) \\ &= \frac{1}{D} [R_i^k \lambda_i^k \lambda_i^{\bar{k}} + R_i^k \lambda_i^k \Delta \lambda_i^{\bar{k}} \\ &\quad + R_i^k \lambda_i^{\bar{k}} \Delta \lambda_i^k + R_i^k \Delta \lambda_i^k \Delta \lambda_i^{\bar{k}}] \end{aligned}$$

Now the first term  $R_i^k \lambda_i^k \lambda_i^{\bar{k}} = \lambda_i^k$ ,

and the third  $R_i^k \lambda_i^{\bar{k}} \Delta \lambda_i^k = \Delta \lambda_i^k$  from (56). For small increments the last term is of the second order

and may be neglected so that

$$u^h = \frac{1}{D} [\lambda^h + \Delta \lambda^h + R \frac{h}{2} \lambda^h \Delta \lambda^{\bar{h}}]$$

To evaluate the constant D, we first sum on h.

$$\begin{aligned} u^h &= \frac{1}{D} [\lambda^h + \Delta \lambda^h + R \frac{h}{2} \lambda^h \Delta \lambda^{\bar{h}}] \\ 1 &= \frac{1}{D} [1 + 0 + \Delta \lambda^{\bar{h}}] = \frac{1}{D} \end{aligned}$$

thus

$$u^h = \lambda^h + \Delta \lambda^h + R \frac{h}{2} \lambda^h \Delta \lambda^{\bar{h}}$$

We see that the offspring of the displaced population is equal to this population plus an additional increment  $R \frac{h}{2} \lambda^h \Delta \lambda^{\bar{h}}$ . Clearly a necessary condition for stability is that this be opposite in sign to the original increment  $\Delta \lambda^h$ . Now the original increment was completely arbitrary. Let us fix h and s as two constant indices and suppose the components of the increment were  $\Delta \lambda^h = +\epsilon$ ,  $\Delta \lambda^s = -\epsilon$  and all other components zero.

Taking  $\epsilon$  positive we have the condition

$$R \frac{h}{2} \lambda^h \Delta \lambda^{\bar{h}} < 0$$

or

$$(R^h \Delta \lambda^h + R^s \Delta \lambda^s + \dots + R^n \Delta \lambda^n) \lambda^h < 0$$

For a realizable population  $\lambda^h$  is positive and all the terms in the parentheses are zero except  $\Delta \lambda^h$  and  $\Delta \lambda^s$  giving

$$R^h \epsilon - R^s \epsilon < 0$$

or

$$R_s^h < R_h^h$$

This proves the theorem.

For a dimorphic factor this condition is also easily shown to be sufficient for stability; but examples show that it is not always sufficient for more than two alleles. Sufficient conditions (not necessary) for any number of alleles are that

$$R_s^h = K_1, \quad , \text{a constant independent of } h \quad (57)$$

$$R_h^h = K_2, \quad , \text{a constant independent of } h \text{ and } s$$

$$K_1 < K_2$$

For then the correction term is

$$\begin{aligned} & \lambda_s^h [ R_s^h \sum_{\substack{i=1 \\ i \neq h}}^n \Delta \lambda_i^h + R_h^h \Delta \lambda_h^h ] \\ &= \lambda_s^h [ -R_s^h \Delta \lambda_h^h + R_h^h \Delta \lambda_h^h ] \\ &= \lambda_s^h \Delta \lambda_h^h (K_1 - K_2) \end{aligned}$$

which is clearly opposite in sign and less in absolute value than  $\Delta \lambda_h^h$ .

## VI. A Calculus of Populations

Up to the present, all our population symbols have been constants, i.e. each represented a certain particular population. The manipulation of these discrete sets of numbers constitutes an algebra. Sometimes, however, it is convenient to consider continuous time variations of a population. Such a study leads to a calculus of populations. We have already used, in the preceding section, the idea of an incremental population.

In this section we will define the "derivative" of a population and develop some of the fundamentals of the calculus.

First let us generalize our idea of population to include variable populations, i.e. populations that are functions of time. We indicate this functional dependence by the usual notation e.g.

$$\lambda_{j^k}^{hi}(t) \quad (58)$$

represents the genetic constitution of the population  $\lambda_{j^k}^{hi}$  at the time t. In case no ambiguity is introduced we will sometimes omit the argument t, it being understood that

$$\lambda_{j^k}^{hi} = \lambda_{j^k}^{hi}(t)$$

We define the derivative of a population as the indexed symbol whose components are the time derivatives of the components of the population in question. Thus:

$$\frac{d}{dt} \lambda_{j^k}^{l_i}(t) = \lim_{\Delta t \rightarrow 0} \frac{\lambda_{j^k}^{l_i}(t+\Delta t) - \lambda_{j^k}^{l_i}(t)}{\Delta t} \quad (59)$$

We assume the population large enough and the variation of  $\lambda$  smooth enough for the limit to exist in a practical sense.

Note that the derivative of a population is not a population. A population has the property  $\lambda^{..} = 1$  while its derivative has the following property:

Theorem XV.

$$\frac{d}{dt} \lambda^{..} = 0 \quad (60)$$

This is true if we first sum on all indices and then take the derivative or vice versa. Both follow immediately on taking the derivative of (6).

As in ordinary calculus we have simple rules for taking derivatives of sums, products, etc. These are all exactly the same as those of ordinary calculus. We have:

Theorem XVI.

1. If  $\lambda_{j^k}^{l^i}(t)$  is a constant

$$\frac{d}{dt} \lambda_{j^k}^{l^i} = 0 \quad (61)$$

2. If

$$\lambda_{j^k}^{l^i} = R_1 u_{j^k}^{l^i} + R_2 v_{j^k}^{l^i} \quad (62)$$

then

$$\frac{d}{dt} \lambda_{j^k}^{l^i} = R_1 \frac{d}{dt} u_{j^k}^{l^i} + R_2 \frac{d}{dt} v_{j^k}^{l^i}$$

3. If

$$\lambda_{j^k}^{l^i} = u_{j^k}^{l^i} \times v_{j^k}^{l^i} \quad (63)$$

then

$$\frac{d}{dt} \lambda_{j^k}^{l^i} = u_{j^k}^{l^i} \times \frac{d}{dt} v_{j^k}^{l^i} + v_{j^k}^{l^i} \times \frac{d}{dt} u_{j^k}^{l^i}$$

The first two of these rules for differentiation are obvious, since by fixing the indices they merely state the ordinary rules for differentiating constants and sums. The third, which is the analogue of Leibnitz rule for differentiating a product, requires proof.

Starting with the definition of a derivative we have:

$$\begin{aligned} \frac{d}{dt} \lambda_{j^k}^{l^i} &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} [\lambda_{j^k}^{l^i}(t + \Delta t) - \lambda_{j^k}^{l^i}(t)] \\ &= \lim_{\Delta t \rightarrow 0} \frac{1}{\Delta t} [ \{u_{j^k}^{l^i}(t) + \Delta u_{j^k}^{l^i}(t)\} \cdot \{v_{j^k}^{l^i}(t) + \Delta v_{j^k}^{l^i}(t)\} - u_{j^k}^{l^i}(t) \cdot v_{j^k}^{l^i}(t) ] \end{aligned}$$

where  $\Delta u_{j^k}^{l_i}(t) = u_{j^k}^{l_i}(t+\Delta t) - u_{j^k}^{l_i}(t)$  and similarly for

$\Delta v_{j^k}^{l_i}(t)$ . Now the first cross product may be

multiplied out by our distributive law (Theorem VIII) giving

$$\lim_{\Delta t \rightarrow 0} [u_{j^k}^{l_i}(t) \times \frac{\Delta v_{j^k}^{l_i}(t)}{\Delta t} + v_{j^k}^{l_i}(t) \frac{\Delta u_{j^k}^{l_i}(t)}{\Delta t} + \Delta u_{j^k}^{l_i}(t) \times \frac{\Delta v_{j^k}^{l_i}(t)}{\Delta t}]$$

The third term in general tends to zero with  $\Delta t$  so that our limit is:

$$u_{j^k}^{l_i} \frac{d}{dt} v_{j^k}^{l_i} + v_{j^k}^{l_i} \frac{d}{dt} u_{j^k}^{l_i}$$

the desired result.

For a population  $\lambda_{j^k}^{l_i}$  intermixing at random this reduces to

$$\frac{d}{dt} (\lambda_{j^k}^{l_i} \times \lambda_{j^k}^{l_i}) = 2 \lambda_{j^k}^{l_i} \times \frac{d}{dt} \lambda_{j^k}^{l_i}$$

A population whose components are analytic functions of time may be expanded in a Taylor series

$$\lambda_{j^k}^{l_i}(t) = \lambda_{j^k}^{l_i}(0) + \frac{d}{dt} \lambda_{j^k}^{l_i} \Big|_{t=0} \cdot t + \frac{d^2}{dt^2} \lambda_{j^k}^{l_i} \Big|_{t=0} \frac{t^2}{2} + \dots$$

for by fixing the indices we are again merely stating the standard Taylor Theorem.

### Bibliography

Inasmuch as no work has been done previously along the specific algebraic lines indicated in this thesis, our references must be of a fairly general nature. For a good introductory treatment of genetics, we recommend:

Sinnott, E.W. and Dunn, L.C. Principles of Genetics.  
McGraw Hill, 1932.

For mathematical treatment of certain phases of hereditary phenomena, the following may be consulted:

Fisher, R.A. The genetical theory of natural selection.  
Oxford Clarendon Press, 1930, Pp. xiv 272.

Wright, Sewall. Genetics, 1921, V.6, Systems of Mating.  
II. The effects of inbreeding on the genetic composition of a population.

Ibid. Pp.162-166. IV. The effects of selection.

Haldane, J.B.S. The cause of evolution. London,  
Longmans Green, 1932.

Robbins, R.B. Some applications of mathematics to breeding problems. I. Genetics, 1917, 2, 489-504;  
II. Ibid, 1918, V.3, 73-92; III. Ibid, 1918, V.3, 375-389.

Hogben, L. A matrix notation for Mendelian populations.  
Proc. Royal Soc. Edinburgh, 1933, 53, 7-25.

### Biography of the Author

Claude E. Shannon was born in Petoskey, Michigan, April 30, 1916. The first sixteen years of his life were spent in Gaylord, Michigan, and there he attended the Public School, graduating from the Gaylord High School in 1932. He entered the University of Michigan in the fall of that year and took an engineering course, specializing in Electrical Engineering and Mathematics. While a senior, he was elected a member of Phi Kappa Phi and an associate member of Sigma Xi. In 1936, he obtained the degrees of Bachelor of Science in Electrical Engineering and Bachelor of Science in Mathematics from the University of Michigan and in the same year started working as a Research Assistant in Electrical Engineering at the Massachusetts Institute of Technology.

From September, 1936 to June, 1938, he was in charge of the Differential Analyzer at the Massachusetts Institute of Technology, a machine used for obtaining numerical solutions to differential equations. During the summer of 1938, he did research work on the design of the Bush Rapid Selector, chiefly connected with the vacuum tube circuits employed in this device. In September, he left the Electrical Engineering Department to become an Assistant in the Department of Mathematics. In the spring of 1939, he was elected to full membership of Sigma Xi. The summer of 1939

was spent at the biological station of the Carnegie Institution of Washington at Cold Spring Harbor, New York, in the study of the application of algebraic methods to hereditary phenomena. In the autumn, he returned to the Massachusetts Institute of Technology as Bolles Fellow and continued work toward a degree of Doctor of Philosophy in Mathematics. In January, 1940, he was married. In the same month, he was awarded the Alfred Noble Prize, of the combined engineering societies of the United States, an award given each year to a person not over thirty for a paper published in one of the journals of the participating societies. This was given in recognition of the paper "A Symbolic Analysis of Relay and Switching Circuits". That spring, he was elected to membership in the American Mathematical Society as an institutional nominee of the Massachusetts Institute of Technology. He is also an associate member of the American Institute of Electrical Engineers and of the Institute of Radio Engineers. He was recently awarded a National Research Fellowship in Mathematics for the year 1940-41.

Published Papers:

"A Symbolic Analysis of Relay and Switching Circuits", Transactions of the American Institute of Electrical Engineers, 1938. Pp. 713-23

**ABSTRACT**

**of**

**"AN ALGEBRA FOR THEORETICAL GENETICS"**

**by**

**Claude E. Shannon**

A Thesis Submitted in Partial Fulfillment  
of the Requirements for the Degree of  
Doctor of Philosophy

**from the**

**Massachusetts Institute of Technology**

**1940**

In this thesis, an algebra is constructed for studying the dynamics of Mendelian populations. The symbols of the algebra represent groups of individuals or populations. The indexed symbol  $\lambda_{j,k}^{k,i}$ , for example, represents a population in which two gene loci are under consideration (the number of loci corresponds to the number of pairs of indices). The number of allelomorphs in each locus is completely arbitrary, as is the recombination value for the loci. The different components of a population symbol, represented by fixing the indices at specific values, are numbers whose values correspond to the fractions of the population with certain genetic formulae. It is convenient in some cases to consider as populations symbols whose components are negative or even complex. Such symbols cannot, of course, represent an actual group of individuals and are called unrealizable populations, but their use sometimes facilitates the solution of problems.

Addition of two population symbols,  $R\lambda_{j,k}^{k,i} + S\mu_{j,k}^{k,i}$ , results in a third population symbol which is defined in such a way as to represent the population obtained by merely combining the original populations in fractional proportions corresponding to the scalar coefficients of R and S. Cross multiplication of population symbols  $\lambda_{j,k}^{k,i} \times \mu_{j,k}^{k,i}$  gives a population symbol which

is defined in such a way as to represent the expected offspring population when the two original populations are crossmated at random. When two gene loci are considered, this is realized by the mathematical definition

$$\lambda_{j^k}^{h_i} \times \mu_{j^k}^{h_i} = \frac{1}{2} [p_0 \lambda_{..}^{h_i} + p_1 \lambda_{..}^{h_i}] [p_0 \mu_{..}^{h_k} + p_1 \mu_{..}^{h_k}] \\ + \frac{1}{2} [p_0 \lambda_{..}^{h_k} + p_1 \lambda_{..}^{h_k}] [p_0 \mu_{..}^{h_i} + p_1 \mu_{..}^{h_i}]$$

in which  $p_1 = 1 - p_0$  is the recombination value for the two loci, and replacing an index by a dot indicates summation of the population symbol on that index. Cross multiplication is defined analogously for  $n$  loci. It is shown that this algebra is commutative on addition and multiplication, distributive, and associative on addition but not on multiplication. These laws together with two fundamental manipulation theorems: one, that summation of a population on all indices gives unity and two, that inverting the upper and lower rows of indices of a population leaves it unchanged, form the basic algorithms of the algebra.

A number of the well known theorems of theoretical genetics are easily proved by means of this algebra. In addition, a number of new results are found. Completely general formulae for the  $n$ th generation offspring under random intermatting of an arbitrary initial population are developed both for the cases of two and of three linked

factors. For two linked factors, the formula for the  $n$ th generation is

$$\lambda_{j,k}^{n+1} = [p_0 \lambda_{..}^{n+1} + p_1 \lambda_{..}^{n+1}] + (1-p_0^{n+1}) \lambda_{..}^n \lambda_{..}^n$$

$$[p_0^n (p_0 \lambda_{..}^{n+1} + p_1 \lambda_{..}^{n+1}) + (1-p_0^n) \lambda_{..}^n \lambda_{..}^n]$$

in which  $\lambda_{j,k}^{n+1}$  is the initial population and  $p_1$  the recombination value. Incidental to this, it is shown that a recombination value  $> \frac{1}{2}$  is impossible when there is no interference. Conditions are found for the stability under random intermating of a population when one or more loci are considered. For the case of one locus, three sets of equivalent necessary and sufficient conditions are established.

By means of certain homogeneous  $\gamma$  populations an arbitrary population may be expanded in a finite series displaying its various components. This expansion, together with the multiplication law for the  $\gamma$  populations, displays the elements of this algebra as hypercomplex numbers. It is shown that an arbitrary population may be expanded uniquely as a sum of any  $n$  linearly independent populations where  $n$  is the number of different possible genetic formulae for the factors considered.

It is possible to write down various types of equations involving known and unknown populations using the

operations of addition and cross multiplication. In general, such an equation can be interpreted as a breeding experiment involving one or more unknown populations and resulting in a genetically known population. Methods are developed whereby most such equations can be solved in case a solution exists. Briefly this method of solution may be summarized as follows. By summing on one or more indices of the unknown populations, enough data about them is obtained to uniquely determine any cross products in which they appear. The cross product terms in the original equations then become known and the equations may be solved in exactly the same way as ordinary linear algebraic equations.

In case a selective action exists favoring individuals of a certain genetic constitution, the previous formulae for stability no longer hold. Although this more difficult problem has not been completely solved, a necessary condition for the possible existence of a stable population under an arbitrary given selective action is established, and a formula for this population is developed. This has only been done for the case of a single locus.

A start has been made toward the development of a calculus of populations, i.e. the study of populations which may vary continuously with time. The time derivative of a population is defined. The derivative of a population, although an indexed symbol, is not itself a popu-

lation. All the ordinary rules of derivation including the Leibnitz rule for the derivative of a cross product of populations are shown to hold true. Also, a population may be expanded in a Taylor series in powers of time, of the same form as the ordinary Taylor series.