

# Prediction of Double Crosses from Single Crosses<sup>\*1</sup>

C. CLARK COCKERHAM

Department of Experimental Statistics, North Carolina State University at Raleigh

**Summary.** A unified theory, which takes into account both the genetic and experimental conditions, for predicting double cross hybrids from single cross hybrids is developed. The method is analogous to that for selection indices. Relationship of the prediction model to the genetic model is explored. JENKINS' (1934) three single cross predictors, a best single cross predictor, and selection on the basis of double cross estimates are compared empirically for an additive and dominance genetic model with varying proportions of experimental error variance and different numbers of hybrids. The differences between fixed and random genetic sample approaches to prediction are discussed.

With many inbred lines it is often impractical to test and compare all possible double cross hybrids. The number of single crosses is considerably fewer, however, and it is logical that they can be used to estimate or predict the performance of the double cross. Questions arise then as to the best methods of prediction and the accuracy of the predictions.

genetics and breeding methodology, will be discussed more fully at the end.

## Single Cross Experiment

To furnish a base for familiarity, consider an experiment of  $r$  replications of all possible single crosses from  $p$  inbred lines. This is the common replicated diallel experiment without reciprocals and parent lines, for which the analysis of variance of plot values is that given in Table 1. The components of variance in Table 1 are  $\sigma^2$  for error,  $\sigma_s^2$  for specific combining ability, and  $\sigma_g^2$  for general combining ability.

## Single Cross Information

For each possible double cross the single crosses may be grouped according to their relationship to the double cross. To the double cross  $AB \cdot CD$ , for example, the two parental single crosses,  $AB$  and  $CD$ , bear the same relationship; the four nonparental single crosses,  $AC$ ,  $AD$ ,  $BC$  and  $BD$ , bear the same relationship; the single crosses,  $AE$ ,  $AF$ , ...,  $BE$ ,  $BF$ , ...,  $CE$ ,  $CF$ , ...,  $DE$ ,  $DF$ , ..., involving only one of the lines,  $A$ ,  $B$ ,  $C$  and  $D$ , bear the same relationship; and the single crosses involving none of the four parent lines bear the same relationship, being zero in this latter case. Consequently, for each double cross, the following four means of the single crosses are considered:

$X_p$  = the mean of the two parental single crosses; e. g.,  $(\overline{AB} + \overline{CD})/2$ .

$X_n$  = the mean of the four nonparental single crosses; e. g.,  $(\overline{AC} + \overline{AD} + \overline{BC} + \overline{BD})/4$ .

$X_l$  = the mean of the single crosses from the four parental lines in combination with all other lines; e. g.,  $(\overline{AE} + \overline{AF} + \dots + \overline{BE} + \overline{BF} + \dots + \overline{CE} + \overline{CF} + \dots + \overline{DE} + \overline{DF} + \dots)/4(p-4)$ . (1)

$X_0$  = the mean of the single crosses from the other lines in all possible combinations; e. g.,  $(\overline{EF} + \overline{EG} + \overline{EH} + \dots + \overline{FG} + \overline{FH} + \dots + \overline{GH} + \dots)/(p-4)(p-5)/2$ .

The divisors indicate the number of single crosses in each mean, and the bar over each single cross indicates that it is a mean for the experiment ( $r$  replications in the example). These means are dependent in that

$$\overline{X} = [2 X_p + 4 X_n + 4(p-4) X_l + (p-4)(p-5) X_0]/p(p-1)/2 \quad (2)$$

Table 1. Analysis of variance of all possible single crosses.

Source	Degrees of Freedom	Mean Square	Expectation of mean square
Replications	$r - 1$		
General	$p - 1$	$\hat{M}_3$	$M_3 = \sigma^2 + r \sigma_s^2 + r(p-2) \sigma_g^2$
Specific	$p(p-3)/2$	$\hat{M}_2$	$M_2 = \sigma^2 + r \sigma_s^2$
Error	$(r-1)(p-2)(p+1)/2$	$\hat{M}_1$	$M_1 = \sigma^2$

The logic of predicting double cross hybrids from single cross hybrids was recognized early in the development of double cross hybrids in corn. JENKINS (1934) presented four alternative methods of predicting double cross performance, three of which involved single crosses, and the first favorable evidence for their use. Since then several studies of one or more of the methods have been reported (reviewed by SPRAGUE, 1955).

EBERHART (1964) gave theoretical relations among hybrids and considered two of JENKINS' (1934) predictors in conjunction with information on three-way crosses in the prediction of double cross performance. EBERHART et al. (1964) compared predicted values of double crosses in maize for various predictors, but did not have double cross information to evaluate the methods. EBERHART's procedures were based on what shall be referred to as a fixed sample approach. In contrast, the theory in the present study, which also takes into account both the genetic and experimental conditions, is based on a random sample approach. The two approaches, both of which have been used extensively in quantitative

\* Dedicated to Dr. GEORGE F. SPRAGUE on the occasion of his 65th birthday.

<sup>1</sup> Contribution from the Department of Experimental Statistics, N. C. Agricultural Experiment Station, Raleigh, North Carolina. Published with the approval of the Director of Research as Paper No. 2330 of the Journal Series. This investigation was supported in part by Public Health Service Research Grant GM 11546 from the Division of General Medical Sciences.

is the mean of all the single crosses and is the same for each double cross. Thus, only three not wholly dependent kinds of single cross information are available for a given double cross.

It is helpful to define two additional means of the single crosses for each double cross.

$$\left. \begin{aligned} X_{\bar{p}} &= (X_p + 2 X_n)/3 = \text{the mean of the six} \\ &\quad \text{single crosses involving the four paren-} \\ &\quad \text{tal lines of the double cross.} \\ X_{\bar{l}} &= [X_p + 2 X_n + (p-4) X_l]/(p-1) \\ &= \text{the mean of the four parental line} \\ &\quad \text{means where a parental line mean is} \\ &\quad \text{the mean of all } (p-1) \text{ single crosses} \\ &\quad \text{involving that line.} \end{aligned} \right\} \quad (3)$$

### Linear Prediction Equation

The result desired is the linear combination of the three kinds of information,  $Z$ 's, on the single crosses which best predicts the performances,  $Y$ 's, of the double crosses.

$$Y = \mu + \hat{Y} + e = \mu + b_1 Z_1 + b_2 Z_2 + b_3 Z_3 + e. \quad (4)$$

The  $b$ 's are the relative weights to be given to the single cross information and  $e = Y - \mu - \hat{Y}$  is the error of prediction. The mean of the double crosses,  $\mu$ , plays no part in the relative ranking of the double crosses.

Within limits, the choice of the  $Z$ 's is a matter of convenience. In the following form they are uncorrelated,

$$\left. \begin{aligned} Z'_1 &= (Z_1 - X_n - X_{\bar{p}}), \\ Z'_2 &= X_{\bar{p}} - 2 X_l + X_0, \\ Z'_3 &= X_{\bar{l}} - \bar{X} \end{aligned} \right\} \quad (5)$$

and have variances,  $\sigma_{Z'_i}^2 = E(Z_i'^2)$ ,

$$\left. \begin{aligned} \sigma_{Z'_1}^2 &= \left( \frac{\sigma^2}{r} + \sigma_s^2 \right) \frac{1}{12} = \frac{M_2}{r} \frac{1}{12}, \\ \sigma_{Z'_2}^2 &= \left( \frac{\sigma^2}{r} + \sigma_s^2 \right) \frac{(p-1)(p-2)}{6(p-4)(p-5)} \\ &= \frac{M_2}{r} \frac{(p-1)(p-2)}{6(p-4)(p-5)}, \\ \sigma_{Z'_3}^2 &= \left[ \frac{\sigma^2}{r} + \sigma_s^2 + (p-2)\sigma_g^2 \right] \frac{(p-2)(p-4)}{4p(p-1)^2} \\ &= \frac{M_3}{r} \frac{(p-2)(p-4)}{4p(p-1)^2}. \end{aligned} \right\} \quad (6)$$

In matrix notation, the variance covariance matrix is

$$\sigma_{ZZ} = \begin{bmatrix} \sigma_{Z'_1}^2 & 0 & 0 \\ 0 & \sigma_{Z'_2}^2 & 0 \\ 0 & 0 & \sigma_{Z'_3}^2 \end{bmatrix}. \quad (7)$$

Let us write formally, for the moment, the covariances of the information on the single crosses with the double crosses,

$$\sigma_{ZY} = (\sigma_{Z_1Y}, \sigma_{Z_2Y}, \sigma_{Z_3Y}). \quad (8)$$

To be used as a criterion of the best predictor is that set of  $b$ 's which minimizes the variance of the prediction error,  $\sigma_e^2 = E(Y - \mu - \hat{Y})^2$ . The minimization procedure is straightforward, analogous to least squares. The  $b$ 's,

$$\underline{b}^{*'} = (b_1^*, b_2^*, b_3^*), \quad (9)$$

that satisfy this criterion are solutions to the following set of equations,

$$\sigma_{ZZ} \underline{b}^* = \sigma_{ZY}, \quad (10)$$

and

$$b_1^* = \frac{\sigma_{Z_1Y}}{\sigma_{Z_1}^2}, \quad b_2^* = \frac{\sigma_{Z_2Y}}{\sigma_{Z_2}^2}, \quad b_3^* = \frac{\sigma_{Z_3Y}}{\sigma_{Z_3}^2}. \quad (11)$$

For any given predictor,  $\hat{Y}_i = \underline{b}' Z_i$ , let  $\bar{Y}_{is}$  be the mean of the set of double crosses corresponding to a selected set of  $\hat{Y}_i$  with mean  $\bar{Y}_{is}$ , and let  $\bar{Y}$  and  $\bar{\hat{Y}}$  be the corresponding means for the entire sample. Based on linear regression, which to be exact requires normal distribution of  $X$  and  $Y$ ,

$$\left. \begin{aligned} \Delta_i &= \bar{Y}_{is} - \bar{Y} = B_Y \hat{Y}_i (\bar{Y}_{is} - \bar{\hat{Y}}_i) \\ B_Y \hat{Y}_i &= \frac{\text{Cov } Y \hat{Y}_i}{\sigma_{\hat{Y}_i}^2} = \frac{b_i' \sigma_{ZY}}{b_i' \sigma_{ZZ} b_i'} \end{aligned} \right\} \quad (12)$$

where  $B_Y \hat{Y}_i$  is the regression coefficient of  $Y$  on  $\hat{Y}_i$ . For truncation selection one may substitute

$$k \sigma_{\hat{Y}_i} = k \sqrt{b_i' \sigma_{ZZ} b_i} = \bar{\hat{Y}}_{is} - \bar{\hat{Y}}_i \quad (13)$$

into (12), and

$$\Delta_i = k \frac{b_i' \sigma_{ZY}}{\sqrt{b_i' \sigma_{ZZ} b_i}} \quad (14)$$

where  $k$  is a function of the intensity of selection. The above relationships (12, 14) hold for any set of  $b$ 's. For the set,  $\underline{b}^*$ , which minimizes the prediction error, substitution of (10) into (14) leads to

$$\Delta_* = k \sqrt{\underline{b}^{*'} \sigma_{ZZ} \underline{b}^*}, \quad (15)$$

and  $\Delta$  is a maximum.

### Genetic Interpretation

It will be mentioned now, and discussed later, that inherent in the previous definition and solution to the best predictor is the assumption that the single crosses and double crosses are random samples for which the procedure will be used. Under certain assumptions for diploids (COCKERHAM, 1963), mainly random samples from a noninbred population and no linkages, the covariances among relatives can be expressed in the following form,

$$C = \alpha \sigma_A^2 + \delta \sigma_D^2 + \alpha^2 \sigma_{AA}^2 + \alpha \delta \sigma_{AD}^2 + \dots, \quad (16)$$

where the subscripts,  $A$  = additive,  $D$  = dominance, indicate the various components of genetic variance. The coefficients  $\alpha$  and  $\delta$  are given by COCKERHAM

Covariance	$\alpha$	$\delta$	Description
$C_{S_2}$	1	1	— Covariance between single cross relatives with both parents common.
$C_{S_1}$	1/2	0	— Covariance between single cross relatives with one parent common only.
$C_{SD_2}$	1/2	1/4	— Covariance between non-parental single cross and double cross relatives, i. e., $AB$ with $A- \cdot B-$ .
$C_{SD_2}$	1/2	0	— Covariance between parental single cross and double cross relatives, i. e., $AB$ with $AB \cdot -$ . Note: $C_{SD_2} = C_{S_1}$ .
$C_{SD_1}$	1/4	0	— Covariance between single cross and double cross relatives with one parent line common only.

(17)

(1961) for all possible types of single, three-way and double cross relatives. Although a more general solution for the coefficients with partially inbred parental lines has been found (COCKERHAM, 1967), we shall consider the parental lines to be homozygous, and shall reproduce here the coefficients for the hybrid relatives under consideration.

It remains to relate the covariances to the statistical parameters. These are

$$\left. \begin{aligned} \sigma_g^2 &= C_{S1}, \\ \sigma_s^2 &= C_{S2} - 2 C_{S1}, \\ \sigma_{Z,Y} &= (C_{SD2} - C_{SD1})/3, \\ \sigma_{Z,Y} &= 2 \sigma_{Z,Y} + (C_{SD2} - 2 C_{SD1}), \\ \sigma_{Z,Y} &= [3 \sigma_{Z,Y} + (\bar{p} - 2) C_{SD1}] \frac{\bar{p} - 4}{\bar{p}(\bar{p} - 1)}. \end{aligned} \right\} \quad (18)$$

Substitution of the appropriate coefficients,  $\alpha$  and  $\delta$ , in (17) into (16) gives the genetic variance composition of the covariances. Further substitution into (18) and then into (11) leads to expressions for the  $b$ 's in terms of genetic and environmental variances and the experimental dimensions,  $\bar{p}$  and  $r$ . While these expressions are not simple, interrelationships of the genetic and prediction models are illuminating.

#### Relationship of the Prediction to the Genetic Model

First of all, note that  $Z_2$  and  $Z_3$  are constant for the three orders,  $AB \cdot CD$ ,  $AC \cdot BD$  and  $AD \cdot BC$ , of double crosses for any four lines. Consequently, the only information about these comparisons is contained in  $Z_1$  and the relative importance of this information is  $b_1$ . The information in  $Z_2$  is of an interaction form, involving the mean of all combinations of four lines,  $X_{\bar{p}}$ , the mean of the four lines in combination with all other lines,  $X_i$ , and the mean of all other lines in combination,  $X_0$ . A two by two table illustrates the point.

	4 lines	Other lines
4 lines	$X_{\bar{p}}$	$X_i$
Other lines	$X_i$	$X_0$

(19)

It,  $Z_2 = X_{\bar{p}} - 2 X_i + X_0$ , is a measure of the interaction of the four lines with the other lines, and the importance of this general nicking of four lines in predicting the corresponding double cross performance is reflected by  $b_2$ . The relative importance of  $Z_3$  depends on how well a line's performance in all single crosses relates to the double cross performance.

The following alternative expressions of the  $b$ 's will be helpful in genetic clarification,

$$\left. \begin{aligned} b_1^* &= \frac{4(C_{SD2} - C_{SD1})}{M_2/r} = \frac{H_1}{M_2/r}, \\ b_2^* &= \frac{(\bar{p} - 4)(\bar{p} - 5)}{(\bar{p} - 1)(\bar{p} - 2)} \times \\ &\quad \times \left[ \frac{4(C_{SD2} - C_{SD1}) + 6(C_{SD2} - 2 C_{SD1})}{M_2/r} \right] \\ &= \frac{(\bar{p} - 4)(\bar{p} - 5)}{(\bar{p} - 1)(\bar{p} - 2)} \left[ b_1 + \frac{6(C_{SD2} - 2 C_{SD1})}{M_2/r} \right] \\ &= \frac{(\bar{p} - 4)(\bar{p} - 5)}{(\bar{p} - 1)(\bar{p} - 2)} \frac{H_2}{M_2/r}, \\ b_3^* &= \frac{\{(\bar{p} - 1)[8(C_{SD2} - C_{SD1}) + 12(C_{SD2} - 2 C_{SD1})] + 4(\bar{p} - 2) C_{SD1}\}}{(\bar{p} - 2) M_3/r} \\ &= \frac{(\bar{p} - 1)^2 2 b_2 M_2}{(\bar{p} - 4)(\bar{p} - 5) M_3} + \frac{4(\bar{p} - 1) C_{SD1}}{M_3/r} \\ &= \frac{(\bar{p} - 1) H_3}{(\bar{p} - 2) (M_3/r)}. \end{aligned} \right\} \quad (20)$$

The genetic variance compositions of the various covariances of relatives, and functions of them, are listed in Table 2 for ready reference. A glance at the genetic numerators,  $H$ 's, of the  $b$ 's shows that for  $b_1^*$  and  $b_2^*$  to be other than 0 requires some domi-

Table 2. Genetic variance composition of various functions of the covariances of relatives.

Covariance	Genetic Components of Variance									
	$\sigma_A^2$	$\sigma_D^2$	$\sigma_{AA}^2$	$\sigma_{AD}^2$	$\sigma_{DD}^2$	$\sigma_{AAA}^2$	$\sigma_{AAD}^2$	$\sigma_{ADD}^2$	$\sigma_{DDD}^2$	...
$C_{S1} = C_{SD2}$	$\frac{1}{2}$	0	$\frac{1}{4}$	0	0	$\frac{1}{8}$	0	0	0	...
$C_{S2}$	1	1	1	1	1	1	1	1	1	...
$C_{SD1}$	$\frac{1}{4}$	0	$\frac{1}{16}$	0	0	$\frac{1}{64}$	0	0	0	...
$C_{SD2}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	...
$C_{S2} - 2 C_{S1}$	0	1	$\frac{1}{2}$	1	1	$\frac{3}{4}$	1	1	1	...
$C_{SD2} - 2 C_{SD1}$	0	0	$\frac{1}{8}$	0	0	$\frac{3}{32}$	0	0	0	...
$C_{SD2} - C_{SD1}$	0	$\frac{1}{4}$	0	$\frac{1}{8}$	$\frac{1}{16}$	0	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{64}$	...
$H_1$	0	1	0	$\frac{1}{2}$	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	...
$H_2$	0	1	$\frac{3}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{9}{16}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	...
$H_3$	$(\bar{p} - 2)$	2	$\frac{\bar{p} + 4}{4}$	1	$\frac{1}{2}$	$\frac{\bar{p} + 16}{16}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	...

nance and/or epistasis. Additive types of epistasis contribute to  $H_2$  but not to  $H_1$ . The dominating coefficients in  $H_3$  are for the all-additive types of variance for large  $\bar{p}$  but all variances are included to some degree.

Since none of the differences among the covariances in (20) can be negative, the following relationships must hold,

$$\left. \begin{aligned} 0 \leq b_1^* \leq \frac{1}{2}, \quad b_1^* \leq \frac{(\bar{p} - 1)(\bar{p} - 2)}{(\bar{p} - 4)(\bar{p} - 5)} b_2^* \leq \frac{3}{2}, \\ b_2^* \frac{12(\bar{p} - 1) M_2}{(\bar{p} - 2) M_3} \leq b_3^* \leq \frac{2(\bar{p} - 1)}{\bar{p} - 2}. \end{aligned} \right\} \quad (21)$$

While these relationships do not definitely establish an order among the  $b$ 's, they have implications for practical applications.

Simplification in the genetic model leads to simplification in the prediction model. If all of the genetic variance is additive variance, the  $b$ 's reduce to

$$\underline{b}^1 = \left[ 0, 0, \frac{(p-1)\sigma_A^2}{(\sigma^2/r) + (p-2)\sigma_A^2} \right]. \quad (22)$$

Since only proportionality of the  $b$ 's are important,

$$\underline{b}^1 \propto (0, 0, 1), \quad (23)$$

and the best prediction model involves only the parental line means,

$$\hat{Y}_1 = X_i - \bar{X}, \quad (24)$$

which is one of the predictors proposed by JENKINS (1934) (his method C). Thus, this predictor is best only when there are only additive effects of genes.

Another of JENKINS' predictors (his method B) is best when all of the genetic variance is dominance variance,

$$\underline{b}^2 = \left[ \frac{\sigma_D^2}{(\sigma^2/r) + \sigma_D^2}, \frac{(p-4)(p-5)}{(p-1)(p-2)} \frac{\sigma_D^2}{((\sigma^2/r) + \sigma_D^2)}, \frac{2(p-1)}{(p-2)} \frac{\sigma_D^2}{((\sigma^2/r) + \sigma_D^2)} \right], \quad (25)$$

or

$$\underline{b}^2 \propto \left[ 1, \frac{(p-4)(p-5)}{(p-1)(p-2)}, \frac{2(p-1)}{p-2} \right], \quad (26)$$

and the prediction equation is

$$\hat{Y}_2 = X_n - \bar{X}, \quad (27)$$

which involves only the nonparental single crosses. A comment is in order. The definitions of genetic effects and variances are such that for there to be only dominance variance requires overdominance at all loci, and further, the noninbred population must have equilibrium frequencies at each locus, a very unlikely condition. There is another case, however, of more importance for which  $\underline{b}^2$  as in (26) is equivalent to  $\underline{b}^*$ . This is when all of the genetic variance is additive and dominance and the experimental error variance is zero, *i. e.*,  $\sigma^2/r = 0$ . Thus, with increasing  $r$ ,  $\underline{b}^2$  tends toward the best predictor when there are both additive and dominance effects of genes but no epistatic effects.

JENKINS' third method (A) of utilizing single crosses assumes the following set of  $b$ 's,

$$\underline{b}^3 \propto \left( 0, \frac{(p-4)(p-5)}{(p-1)(p-2)}, \frac{2(p-1)}{p-2} \right), \quad (28)$$

which corresponds to the following prediction equation,

$$\hat{Y}_3 = X_{\bar{p}} - \bar{X} \quad (29)$$

and amounts to using the mean of the six single crosses of the four parents of the double cross. There appears to be no simple genetic situation for which this predictor would be best.

If there are no epistatic effects, but allowing for both additive and dominance variance,

$$\underline{b}^4 = \left\{ \frac{\sigma_D^2}{(\sigma^2/r) + \sigma_D^2}, \frac{(p-4)(p-5)}{(p-1)(p-2)} \frac{\sigma_D^2}{((\sigma^2/r) + \sigma_D^2)}, \frac{(p-1)[2\sigma_D^2 + (p-2)\sigma_A^2]}{(p-2)[(\sigma^2/r) + \sigma_D^2 + (p-2)\sigma_A^2/2]} \right\}, \quad (30)$$

or

$$\underline{b}^4 \propto \left\{ 1, \frac{(p-4)(p-5)}{(p-1)(p-2)}, \frac{(p-1)[(\sigma^2/r) + \sigma_D^2][2\sigma_D^2 + (p-2)\sigma_A^2]}{(p-2)\sigma_A^2[(\sigma^2/r) + \sigma_D^2 + (p-2)\sigma_A^2/2]} \right\}. \quad (31)$$

All of the parameters are contained in the diallel experiment. Consequently, the  $b$ 's may be written as functions of the expected mean squares,

$$\underline{b}^4 = \left[ \frac{M_2 - M_1}{M_2}, \frac{(p-4)(p-5)(M_2 - M_1)}{(p-1)(p-2)M_2}, \frac{(p-1)(M_3 + M_2 - 2M_1)}{(p-2)M_3} \right], \quad (32)$$

or

$$\underline{b}^4 \propto \left[ 1, \frac{(p-4)(p-5)}{(p-1)(p-2)}, \frac{2(p-1)M_2}{(p-2)M_3} \left( 1 + \frac{M_3 - M_2}{M_2 - M_1} \right) \right]. \quad (33)$$

Thus, from a diallel experiment, assuming no epistasis, the  $b$ 's may be estimated. The proportional  $b$ 's in (31, 33) assume at least some dominance variance. Otherwise we have divided through by zero in standardizing  $b_1$  to be one.

To admit all types of epistatic variance requires estimates not available in the diallel experiment, namely,  $H_1$ ,  $H_2$  and  $H_3$  or two relative values, say  $H_2/H_1$  and  $H_3/H_1$ . If there are no additive types of epistatic variance then  $H_2 = H_1$  and  $H_3 = 2H_1 + 2(M_3 - M_2)/r$ , but still required is the relative value of  $\sigma_A^2 = 2(M_3 - M_2)/(p-2)r$  to  $H_1$ . The proportional  $b$ 's would be the same as in (33) except for  $b_3$ ,

$$\underline{b}^5 \propto \left[ 1, \frac{(p-4)(p-5)}{(p-1)(p-2)}, \frac{2(p-1)M_2}{(p-2)M_3} \left( 1 + \frac{M_3 - M_2}{rH_1} \right) \right], \quad (34)$$

and  $b_3$  would be relatively larger when there was epistasis involving dominance instead of only dominance effects.

### Goodness of Predictors

The appeal of best may not always be warranted. Sometimes a best predictor as defined herein is only infinitesimally better than another one, and all may be good or poor depending on the parametric situation.

We shall compare predictors empirically as ratios of their expected gains,  $\Delta$ 's, to gain expected from selecting among the true values,  $Y$ 's, of the double crosses. The gain expected from selecting among the  $Y$ 's,

$$\Delta_Y = k\sigma_Y, \quad (35)$$

is dependent upon the selection intensity reflected in  $k$  and  $\sigma_Y^2$  which is derived in an appendix. The efficiency measure,

$$\frac{\Delta_i}{\Delta_Y} = \frac{b'_i \sigma_{ZY}}{\sqrt{b'^i \sigma_{ZZ} b^i \sigma_Y}} = \rho_{iY}, \quad (36)$$

is also the correlation,  $\rho_{iY}$ , between the predictor,  $\hat{Y}_i$ , and  $Y$ , which can be verified from relationships given in (12). The relative efficiency of any two predictors,  $\hat{Y}_i$  and  $\hat{Y}_j$ , may be obtained simply as

$$\frac{\Delta_i}{\Delta_j} = \frac{\rho_{iY}}{\rho_{jY}} \quad (37)$$

which, when  $j = *$ ,

$$\frac{\Delta_i}{\Delta_*} = \frac{\rho_{iY}}{\rho_{*Y}} = \rho_{i*}, \quad (38)$$

is also the correlation between  $\hat{Y}_i$  and  $\hat{Y}_*$ .

Since the true genotypic values of the double crosses are never known, we wish to know also how

effective is selection based on experimental estimates of double cross performance as a further basis of comparison. Among many alternatives we shall consider that the same amount of experimental information is available for the double crosses as for the single crosses, *i. e.*, the same number,  $r$ , of replicate plots for each double cross. Denoting the mean over replications of individual double crosses as  $\tilde{Y}$ ,

$$\tilde{Y} = Y + \bar{\varepsilon}, \quad (39)$$

where  $\bar{\varepsilon}$  is the experimental error. The gain from selecting on  $\tilde{Y}$ , following arguments (12) through (14), is

$$\Delta_{\tilde{Y}} = k \frac{\sigma_{\tilde{Y}}^2}{\sigma_{\tilde{Y}}}. \quad (40)$$

The variance,  $\sigma_{\tilde{Y}}^2$ , among the estimates is derived simultaneously with  $\sigma_{\tilde{Y}}^2$  in the Appendix.

The ratio of the gain (40) from using the estimates to that (35) from knowing the true values,

$$\frac{\Delta_{\tilde{Y}}}{\Delta_Y} = \frac{\sigma_Y}{\sigma_{\tilde{Y}}} = \rho_{\tilde{Y}Y}, \quad (41)$$

is again the correlation between the two, which approaches one as the experimental error variance approaches zero.

To be compared empirically are JENKINS' (1934) three predictors (23, 26 and 28) and the best predictor,  $\hat{Y}_*$ . Their efficiencies in the form of (36) and the efficiency of experimental double cross prediction (41) give five comparative measures,

$$\rho_{1Y}, \rho_{2Y}, \rho_{3Y}, \rho_{*Y}, \rho_{\tilde{Y}Y}. \quad (42)$$

These efficiency measures are affected only by changes in the relative values of the variance parameters. Consequently, if there are no epistatic effects, the variances can be reduced to two quantities,

$$\beta = \frac{\sigma_A^2 + \sigma_D^2}{(\sigma^2/r) + \sigma_A^2 + \sigma_D^2}, \quad \gamma = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_D^2}, \quad (43)$$

such that

$$\left[ \frac{\sigma^2}{r} = 1 - \beta \right] + [\sigma_D^2 = \beta (1 - \gamma)] + [\sigma_A^2 = \gamma \beta] = 1. \quad (44)$$

The proportion of the total variance that is genotypic is measured by  $\beta$ , and of the genotypic variance that is additive by  $\gamma$ .

The five efficiency measures (42) are given for various combinations of  $\gamma$  and  $\beta$  in Table 3 for  $p = 6$  and  $p = 20$ , and in Table 4 for  $p = 100$  and  $p = \infty$ . Also included in the Tables is  $\sigma_Y$  to show how it, and consequently  $\Delta_Y$  in (35), varies with the parameters. While the error variance,  $\sigma_{\varepsilon}^2$ , in  $\sigma_{\tilde{Y}}^2$  for double crosses is slightly larger than  $\sigma^2$  for the single crosses, (74) and (75), the two were taken to be the same which has a very slight inflationary effect on  $\rho_{\tilde{Y}Y}$ .

First, consider the  $\rho_{iY}$ 's in Tables 3 and 4 relating to single cross prediction only. As a base of reference,  $\hat{Y}_*$  is an ideal and represents the best that one can do with single cross prediction. Its efficiency,  $\rho_{*Y}$ , increases with  $\gamma$ , with  $\beta$ , and with  $p$  except for  $\gamma = 0$ . There is an overriding effect of  $\beta$ , however, in that variations in  $\gamma$  have less of an effect as  $\beta$

Table 3. *Efficiencies of Predictors for  $p = 6$  and  $p = 20$  (in parentheses).*

		$\gamma$							
		0.0	0.1	0.3	0.5	0.7	0.9	1.0	
$\beta$	0.1	$\rho_{1Y}$ .21 (.13)	.23 (.23)	.28 (.38)	.32 (.50)	.36 (.59)	.40 (.67)	.43 (.71)	
		$\rho_{2Y}$ .32 (.32)	.31 (.33)	.31 (.34)	.31 (.36)	.31 (.37)	.30 (.39)	.30 (.39)	
		$\rho_{3Y}$ .23 (.26)	.25 (.28)	.28 (.33)	.31 (.38)	.35 (.41)	.38 (.45)	.40 (.46)	
		$\rho_{*Y}$ .32 (.32)	.32 (.34)	.32 (.42)	.34 (.51)	.37 (.60)	.41 (.67)	.43 (.71)	
		$\rho_{\tilde{Y}Y}$ .14 (.16)	.14 (.17)	.14 (.18)	.14 (.19)	.14 (.19)	.14 (.20)	.14 (.21)	
		$\sigma_Y$ .14 (.16)	.13 (.16)	.13 (.17)	.13 (.18)	.13 (.19)	.13 (.20)	.13 (.20)	
	0.3	$\rho_{1Y}$ .37 (.23)	.40 (.37)	.47 (.55)	.53 (.69)	.59 (.78)	.65 (.86)	.68 (.89)	
		$\rho_{2Y}$ .55 (.55)	.55 (.56)	.54 (.58)	.54 (.60)	.54 (.62)	.53 (.63)	.53 (.64)	
		$\rho_{3Y}$ .40 (.44)	.43 (.48)	.48 (.55)	.53 (.61)	.57 (.66)	.62 (.70)	.65 (.72)	
		$\rho_{*Y}$ .55 (.55)	.55 (.57)	.55 (.65)	.57 (.72)	.61 (.80)	.65 (.86)	.68 (.89)	
		$\rho_{\tilde{Y}Y}$ .27 (.31)	.27 (.32)	.27 (.33)	.27 (.35)	.26 (.36)	.26 (.38)	.26 (.38)	
		$\sigma_Y$ .23 (.27)	.23 (.28)	.23 (.30)	.23 (.31)	.23 (.33)	.22 (.34)	.22 (.35)	
	0.5	$\rho_{1Y}$ .48 (.30)	.51 (.44)	.59 (.63)	.65 (.76)	.72 (.85)	.78 (.92)	.82 (.95)	
		$\rho_{2Y}$ .71 (.71)	.71 (.72)	.70 (.74)	.70 (.75)	.70 (.77)	.69 (.78)	.69 (.79)	
		$\rho_{3Y}$ .52 (.57)	.55 (.62)	.61 (.69)	.66 (.74)	.71 (.79)	.76 (.83)	.79 (.84)	
		$\rho_{*Y}$ .71 (.71)	.71 (.73)	.71 (.77)	.72 (.82)	.75 (.87)	.79 (.92)	.82 (.95)	
		$\rho_{\tilde{Y}Y}$ .40 (.44)	.40 (.45)	.39 (.48)	.39 (.49)	.39 (.51)	.38 (.53)	.38 (.53)	
		$\sigma_Y$ .30 (.35)	.30 (.36)	.30 (.38)	.30 (.40)	.29 (.42)	.29 (.44)	.29 (.45)	
	0.7	$\rho_{1Y}$ .56 (.36)	.60 (.50)	.68 (.68)	.75 (.80)	.81 (.88)	.88 (.95)	.91 (.98)	
		$\rho_{2Y}$ .84 (.84)	.84 (.84)	.83 (.86)	.83 (.87)	.83 (.88)	.83 (.89)	.82 (.89)	
		$\rho_{3Y}$ .62 (.68)	.65 (.72)	.70 (.79)	.76 (.84)	.81 (.88)	.87 (.91)	.89 (.92)	
		$\rho_{*Y}$ .84 (.84)	.84 (.85)	.84 (.87)	.84 (.90)	.86 (.92)	.88 (.95)	.91 (.98)	
		$\rho_{\tilde{Y}Y}$ .55 (.60)	.55 (.61)	.55 (.64)	.54 (.66)	.54 (.67)	.53 (.69)	.53 (.69)	
		$\sigma_Y$ .36 (.41)	.36 (.43)	.35 (.45)	.35 (.48)	.35 (.50)	.34 (.52)	.34 (.53)	
	0.9	$\rho_{1Y}$ .64 (.40)	.68 (.54)	.75 (.71)	.82 (.82)	.88 (.90)	.94 (.97)	.97 (.99)	
		$\rho_{2Y}$ .95 (.95)	.95 (.95)	.95 (.96)	.95 (.96)	.95 (.96)	.94 (.97)	.94 (.97)	
		$\rho_{3Y}$ .70 (.77)	.73 (.81)	.78 (.86)	.84 (.90)	.89 (.94)	.94 (.97)	.97 (.98)	
		$\rho_{*Y}$ .95 (.95)	.95 (.95)	.95 (.96)	.95 (.96)	.95 (.97)	.96 (.98)	.97 (.99)	
		$\rho_{\tilde{Y}Y}$ .79 (.83)	.79 (.84)	.79 (.85)	.79 (.86)	.78 (.87)	.78 (.88)	.78 (.88)	
		$\sigma_Y$ .41 (.47)	.40 (.48)	.40 (.51)	.40 (.54)	.39 (.56)	.39 (.59)	.39 (.60)	
	1.0	$\rho_{1Y}$ .67 (.43)	.71 (.55)	.78 (.72)	.85 (.83)	.91 (.91)	.97 (.97)	1.00 (1.00)	
		$\rho_{2Y}$ 1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	
		$\rho_{3Y}$ .74 (.81)	.77 (.84)	.82 (.89)	.87 (.93)	.92 (.96)	.97 (.99)	1.00 (1.00)	
		$\rho_{*Y}$ 1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	
		$\rho_{\tilde{Y}Y}$ 1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	
		$\sigma_Y$ .43 (.49)	.43 (.51)	.42 (.54)	.42 (.57)	.41 (.59)	.41 (.62)	.41 (.63)	

Table 4. *Efficiencies of Predictors for  $p = 100$  and  $p = \infty$  (in parentheses).*

		$\gamma$						
		0.0	0.1	0.3	0.5	0.7	0.9	1.0
$\beta$	0.1 $\varrho_1 Y$	.06 (.00)	.29 (.43)	.55 (.68)	.70 (.82)	.81 (.91)	.89 (.97)	.92 (1.00)
	$\varrho_2 Y$	.32 (.32)	.33 (.33)	.35 (.36)	.37 (.38)	.39 (.40)	.41 (.42)	.42 (.43)
	$\varrho_3 Y$	.26 (.26)	.29 (.29)	.35 (.35)	.40 (.40)	.44 (.45)	.48 (.48)	.49 (.50)
	$\varrho_* Y$	.32 (.32)	.39 (.51)	.58 (.71)	.71 (.83)	.81 (.91)	.89 (.97)	.92 (1.00)
	$\varrho_{\tilde{Y}} Y$	.16 (.16)	.17 (.17)	.19 (.19)	.20 (.20)	.21 (.21)	.22 (.22)	.23 (.23)
	$\sigma_Y$	.16 (.16)	.17 (.17)	.18 (.18)	.19 (.19)	.20 (.21)	.21 (.22)	.22 (.22)
	0.3 $\varrho_1 Y$	.11 (.00)	.39 (.43)	.64 (.68)	.78 (.82)	.88 (.91)	.95 (.97)	.98 (1.00)
	$\varrho_2 Y$	.55 (.55)	.56 (.57)	.59 (.60)	.62 (.63)	.64 (.65)	.66 (.67)	.67 (.68)
	$\varrho_3 Y$	.45 (.45)	.49 (.50)	.57 (.58)	.63 (.64)	.68 (.69)	.73 (.73)	.74 (.75)
	$\varrho_* Y$	.55 (.55)	.61 (.64)	.73 (.77)	.82 (.85)	.89 (.92)	.95 (.97)	.98 (1.00)
	$\varrho_{\tilde{Y}} Y$	.31 (.31)	.32 (.32)	.35 (.35)	.37 (.37)	.39 (.39)	.40 (.41)	.41 (.42)
	$\sigma_Y$	.27 (.27)	.29 (.29)	.31 (.31)	.33 (.34)	.35 (.36)	.37 (.38)	.38 (.39)
	0.5 $\varrho_1 Y$	.14 (.00)	.42 (.43)	.67 (.68)	.80 (.82)	.90 (.91)	.96 (.97)	.99 (1.00)
	$\varrho_2 Y$	.71 (.71)	.72 (.72)	.75 (.75)	.77 (.77)	.79 (.79)	.80 (.81)	.81 (.82)
	$\varrho_3 Y$	.58 (.58)	.63 (.63)	.71 (.71)	.76 (.77)	.81 (.82)	.85 (.85)	.86 (.87)
	$\varrho_* Y$	.71 (.71)	.74 (.75)	.81 (.83)	.87 (.88)	.92 (.93)	.97 (.98)	.99 (1.00)
	$\varrho_{\tilde{Y}} Y$	.45 (.45)	.46 (.46)	.49 (.50)	.52 (.52)	.54 (.55)	.56 (.57)	.57 (.58)
	$\sigma_Y$	.35 (.35)	.37 (.37)	.40 (.40)	.43 (.43)	.45 (.46)	.48 (.49)	.49 (.50)
	0.7 $\varrho_1 Y$	.17 (.00)	.44 (.43)	.68 (.68)	.81 (.82)	.90 (.91)	.97 (.97)	1.00 (1.00)
	$\varrho_2 Y$	.84 (.84)	.85 (.85)	.87 (.87)	.88 (.88)	.89 (.89)	.90 (.90)	.90 (.91)
	$\varrho_3 Y$	.68 (.68)	.73 (.73)	.80 (.80)	.85 (.86)	.89 (.89)	.92 (.92)	.93 (.94)
	$\varrho_* Y$	.84 (.84)	.85 (.86)	.89 (.89)	.91 (.92)	.94 (.95)	.97 (.98)	1.00 (1.00)
	$\varrho_{\tilde{Y}} Y$	.61 (.61)	.62 (.63)	.65 (.66)	.68 (.68)	.70 (.71)	.72 (.73)	.73 (.73)
	$\sigma_Y$	.42 (.42)	.44 (.44)	.47 (.48)	.51 (.51)	.54 (.55)	.57 (.58)	.58 (.59)
	0.9 $\varrho_1 Y$	.19 (.00)	.45 (.43)	.69 (.68)	.82 (.82)	.91 (.91)	.97 (.97)	1.00 (1.00)
	$\varrho_2 Y$	.95 (.95)	.95 (.95)	.96 (.96)	.96 (.96)	.97 (.97)	.97 (.97)	.97 (.97)
	$\varrho_3 Y$	.77 (.77)	.81 (.82)	.87 (.88)	.92 (.92)	.95 (.95)	.97 (.97)	.98 (.98)
	$\varrho_* Y$	.95 (.95)	.95 (.95)	.96 (.96)	.97 (.97)	.97 (.98)	.98 (.99)	1.00 (1.00)
	$\varrho_{\tilde{Y}} Y$	.83 (.83)	.84 (.84)	.86 (.86)	.88 (.88)	.89 (.89)	.90 (.90)	.90 (.90)
	$\sigma_Y$	.47 (.47)	.50 (.50)	.54 (.54)	.57 (.58)	.61 (.62)	.64 (.65)	.66 (.67)
	1.0 $\varrho_1 Y$	.20 (.00)	.46 (.43)	.69 (.68)	.82 (.82)	.91 (.91)	.97 (.97)	1.00 (1.00)
	$\varrho_2 Y$	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)
	$\varrho_3 Y$	.82 (.82)	.85 (.85)	.90 (.91)	.94 (.94)	.97 (.97)	.99 (.99)	1.00 (1.00)
	$\varrho_* Y$	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)
	$\varrho_{\tilde{Y}} Y$	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)	1.00 (1.00)
	$\sigma_Y$	.50 (.50)	.52 (.52)	.56 (.57)	.60 (.61)	.64 (.65)	.68 (.69)	.69 (.71)

becomes larger. Turning now to  $\hat{Y}_2$  we knew from (26), and the related comments, that it was equivalent to the most efficient predictor,  $\hat{Y}_*$ , when all of the genetic variance was dominance, *i. e.*,  $\gamma = 0$ , and when  $\beta = 1$ , *i. e.*,  $\sigma^2/r = 0$ . Actually, the efficiency is relatively insensitive to changes in  $\gamma$  and  $p$ , and mainly dependent on  $\beta$ . Varying  $p$  has no effect on the efficiency when  $\gamma = 0$ , and there is apparently a  $p$  between 6 and 20 for which variations in  $\gamma$  have an imperceptible, if any, effect on the efficiency because  $\varrho_2 Y$  increases for  $p = 20$  and decreases for  $p = 6$  with an increase in  $\gamma$ .

Considerably more variation is obtained in the efficiencies of  $\hat{Y}_1$ , with all parameters  $\beta$ ,  $\gamma$  and  $p$  having pronounced effects in certain circumstances. This predictor is the most efficient when  $\gamma = 1$ , which we knew from (23). The efficiency increases with  $\gamma$ , increases with  $p$  for large  $\gamma$  but decreases with  $p$  for small  $\gamma$ , and increases with  $\beta$  but  $\beta$  has less of an effect as  $p$  becomes larger and has no effect when  $p = \infty$ . The efficiency is appreciably greater than  $\varrho_2 Y$  only for small  $\beta$  and no smaller than intermediate  $\gamma$ , and for large  $\gamma$  and no larger than intermediate  $\beta$ .

The set of coefficients for  $\hat{Y}_3$  is in some respects intermediate to those for  $\hat{Y}_2$  and  $\hat{Y}_1$ . For most situations in Tables 3 and 4,  $\varrho_3 Y$  is intermediate to  $\varrho_2 Y$  and  $\varrho_1 Y$ . However, near the value of  $\gamma$  for which  $\varrho_1 Y$  and  $\varrho_2 Y$  switch relative magnitudes  $\varrho_3 Y$

is less than both. This feature is demonstrated only for  $p \geq 20$  because of rounding and of the large intervals used for  $\gamma$ .

It should be noted that there is a variety of combinations of  $\gamma$  and  $\beta$  for which all three efficiencies,  $\varrho_1 Y$ ,  $\varrho_2 Y$  and  $\varrho_3 Y$ , are fairly similar. These combinations include intermediate  $\gamma$  and intermediate to small  $\beta$  which may well prevail in many experiments. Also, although none of the three is the best, none of the three efficiencies is far below that of the best.

The efficiency,  $\varrho_{\tilde{Y}} Y$ , of selecting on the basis of double cross estimates, like  $\varrho_2 Y$ , increases with  $\beta$  and  $p$ , and is affected little by variations in  $\gamma$ . Also, as we already knew,  $\varrho_{\tilde{Y}} Y$  approaches one as  $\beta$  approaches one, since when  $\sigma^2/r = 0$  selection is among the true values of the double crosses. The surprising feature is that  $\varrho_{\tilde{Y}} Y$  is never greater than  $\varrho_2 Y$  and for many likely combinations of  $\gamma$  and  $\beta$  it is considerably less. Also, the other predictors,  $\hat{Y}_1$  and  $\hat{Y}_2$ , are more efficient than  $\tilde{Y}$  for most situations.

The standard deviation,  $\sigma_Y$ , among the true values behaves similar to  $\varrho_2 Y$ . It is mainly affected by and increases with  $\beta$ . It increases with  $p$  and much more so for large  $\gamma$  than for small  $\gamma$ . It increases with  $\gamma$  for large  $p$  but the reverse happens for small  $p$ . This interaction of  $p$  and  $\gamma$  stems from a reduction in  $\sigma_Y$  with an increasing proportion of nonadditive variance, but which is offset for small  $p$  by the finite correcting coefficient of the additive variance (see 73

in the Appendix). When  $p = \infty$  the variance is that among unrelated double crosses, i. e.,  $\sigma_Y^2 = \beta(1 + \gamma)/4$ .

### Discussion

A comparison of the fixed and random sample approach is pertinent to the remaining discussion. Let the model for the double crosses be

$$Y = \mu + G, \quad (45)$$

where  $G$  is the genetic effect. Let the model for two alternative single cross predictors be

$$\begin{aligned} X_1 &= \mu_1 + G + e_1, \\ X_2 &= \mu_1 + G + g + e_2. \end{aligned} \quad (46)$$

Provision is made for the experimental errors to be different. One of the predictors,  $X_1$ , is unbiased genetically and has a prediction error of

$$e_1 \quad (47)$$

with variance

$$E(e_1^2) = \sigma_{e_1}^2. \quad (48)$$

Any constants involving the means fall out in comparisons among the  $Y$ 's and their corresponding  $X$ 's, and they will be ignored. The other predictor,  $X_2$ , is biased genetically,  $g$ , and has a prediction error of

$$g + e_2 \quad (49)$$

with an average quadratic value of

$$E(g^2) + \sigma_{e_2}^2. \quad (50)$$

In both cases, random or fixed, the genotypes and environments are assumed uncorrelated.

From a fixed sample standpoint, it is a particular set of single and double crosses to which the procedure is to be applied. Being unable to evaluate each  $g$  in most cases, primary emphasis is placed on finding the least biased predictors with minimum experimental error variance. None of the single cross predictors is unbiased genetically if there is epistasis, as was shown by EBERHART (1964), and only  $\hat{Y}_2$  is unbiased if there are only additive and dominance effects of genes.

Now suppose that the procedure is viewed as being used repeatedly for samples of single and double crosses. Then in the parametric population  $g$  will have a mean of 0 and variance of  $E(g^2) = \sigma_g^2$  even though the parametric population may be a peculiar one consisting only of single and double crosses from highly selected inbred lines. Further,  $G$  will be distributed with a mean of 0 and variance,  $\sigma_G^2$ . For generality, let  $g$  and  $G$  be correlated with covariance,  $\sigma_{gG}$ . With this viewpoint, the prediction error variance is  $\sigma_{e_1}^2$  for  $X_1$  as before and  $\sigma_g^2 + \sigma_{e_2}^2$  for  $X_2$ . Now consider alternative predictors  $c_1 X_1$  and  $c_2 X_2$  where  $c_1$  and  $c_2$  are constants. The constants which minimize the prediction error variances,

$$\left. \begin{aligned} \sigma_{Y \cdot c_1 X_1}^2 &= E[Y - \mu - c_1(X_1 - \mu_1)]^2 \\ &= (1 - c_1)^2 \sigma_G^2 + c_1^2 \sigma_{e_1}^2, \\ \sigma_{Y \cdot c_2 X_2}^2 &= E[Y - \mu - c_2(X_2 - \mu_1)]^2 \\ &= (1 - c_2)^2 \sigma_G^2 - 2c_2(1 - c_2)\sigma_{gG} \\ &\quad + c_2^2(\sigma_g^2 + \sigma_{e_2}^2), \end{aligned} \right\} \quad (51)$$

are the regression coefficients

$$\left. \begin{aligned} c_1^* &= B_{Y X_1} = \frac{\sigma_G^2}{\sigma_G^2 + \sigma_{e_1}^2} = \rho_{G_1} \frac{\sigma_G}{\sigma_{X_1}}, \\ c_2^* &= B_{Y X_2} = \frac{\sigma_G^2 + \sigma_{gG}}{\sigma_G^2 + 2\sigma_{gG} + \sigma_g^2 + \sigma_{e_2}^2} = \rho_{G_2} \frac{\sigma_G}{\sigma_{X_2}} \end{aligned} \right\} \quad (52)$$

and the prediction error variances are

$$\left. \begin{aligned} \sigma_{Y \cdot c_1^* X_1}^2 &= \sigma_G^2(1 - \rho_{G_1}^2), \\ \sigma_{Y \cdot c_2^* X_2}^2 &= \sigma_G^2(1 - \rho_{G_2}^2). \end{aligned} \right\} \quad (53)$$

The gains from selecting on the basis of the predictors, assuming linearity between the  $X$ 's and  $Y$ 's,

$$\left. \begin{aligned} \Delta_{X_1} &= k \frac{\sigma_G^2}{\sqrt{\sigma_G^2 + \sigma_{e_1}^2}} = k c_1^* \sigma_{X_1} = \Delta_{c_1^* X_1} = k \rho_{G_1} \sigma_G, \\ \Delta_{X_2} &= k \frac{\sigma_G^2 + \sigma_{gG}}{\sqrt{\sigma_G^2 + \sigma_{gG} + \sigma_g^2 + \sigma_{e_2}^2}} = k c_2^* \sigma_{X_2} \\ &= \Delta_{c_2^* X_2} = k \rho_{G_2} \sigma_G, \end{aligned} \right\} \quad (54)$$

are invariant of the constant multipliers. The prediction error variances when in their minimum form (53) rank exactly in reverse order to the gains (54). In other than minimum form, the prediction error variances may not rank the predictors the same as the gains, although they do for  $c_1 = c_2 = 1$  when  $\sigma_{gG} = 0$ , which is a simpler situation.

The foregoing illustrates the connection between minimum prediction error variances and gains when applied to random samples. The extension to the simultaneous weighting of two or more kinds of information, for example  $d_1 X_1 + d_2 X_2$ , into a single predictor is straightforward. The constants,  $d_1^*$  and  $d_2^*$ , which maximize the gain also minimize the prediction error variance, although any other set of constants,  $t d_1^*$  and  $t d_2^*$ , where  $t$  is a positive constant, will give the same gain. The maximization of gain is simply an adaptation of the selection index procedure given in considerable detail by HENDERSON (1963), who also pointed out many applications. The same principles are involved in weighting of information from individuals and families in making selections considered by LUSH (1947).

It is seen then that genetically biased predictors may have less prediction error variance than unbiased ones depending upon the bias variance, the covariance between the bias and the true values, and the experimental error variances. These factors are automatically accommodated in the relative weighting of the various single cross means. The number of parental lines,  $p$ , determining the number of single crosses in some of the means plays an important part in the relative weights and in the gains, a reflection of its joint effect on experimental variances and genetic covariances of hybrid means.

The choice among predictors and minimization of prediction error variance in random samples always depends upon a knowledge of the variance parameters, which are never known exactly. For single cross prediction these are  $\sigma_{ZZ}$  in (7) and  $\sigma_{ZY}$  in (8). If one is not willing to assume the hybrids to be members of a random mating population, then  $\sigma_{ZZ}$  and  $\sigma_{ZY}$  must be estimated directly from single crosses and double crosses. If the hybrids are assumed to be members of a random mating population, then estimates of  $\sigma_{ZZ}$  and  $\sigma_{ZY}$  can be constructed from estimates of genetic and environmental variances

from other types of experiments (COCKERHAM, 1963). In any case, estimates are involved. The problem of using estimated variances to construct the  $b$ 's (WILLIAMS, 1962) is that they must be reliable, otherwise the predictor or index may be useless. However, the hybrid prediction framework is much more restricted than that for the general selection index considered by WILLIAMS. The weights,  $b^*$ , cannot be negative and must conform to relative bounds (21). PATEL (1962), in a somewhat similar situation of selecting among single crosses, showed that estimated weights when restricted to fall within known bounds led to efficient selection.

The empirical evaluations were made only for additive and dominance variance and assuming the hybrids to be members of a noninbred population. To compare the predictors empirically with much generality in the epistatic model would be a large task. Epistasis, like dominance, reduces  $\rho_{*Y}$ , and also  $\rho_{\tilde{Y}Y}$  and  $\sigma_Y$  except for small  $p$ , as compared to an additive genetic model. Consequently, the gains,  $\Delta_* \propto \rho_{*Y} \sigma_Y$ ,  $\Delta_{\tilde{Y}} \propto \rho_{\tilde{Y}Y} \sigma_Y$  and  $\Delta_Y \propto \sigma_Y$ , are all reduced by epistasis for a given proportion of the total variance that is environmental. The amount of reduction depends on the kinds and relative amounts of epistatic variance, and will increase with higher order epistasis and more with dominance types than with additive types.

One of the surprising outcomes is the efficiency of all the single cross predictors relative to double cross estimates in Tables 3 and 4. It was pointed out previously that this relative advantage decreases some with epistasis. Also, the most efficient procedure for selecting on the  $\tilde{Y}$ 's has not been used. For each double cross there are eight subdivisions of the double cross estimates (see the Appendix), each related differently to the double cross. The most efficient method is to weight the means of these subdivisions, or correspondingly appropriate linear functions of them, just as was done for the single crosses in predicting double crosses, and in the same way that PATEL (1962) did in selecting on single cross estimates for single crosses. The most efficient double cross predictor should always be superior to single cross prediction. In practice, however, one is concerned with the prediction of double crosses without having to obtain  $p(p-1)(p-2)(p-3)/8$  estimates.

Whether to choose one of JENKINS' (1934) three predictors,  $\hat{Y}_1$ ,  $\hat{Y}_2$  and  $\hat{Y}_3$ , or to estimate the best one,  $\hat{Y}_*$ , is probably unimportant in practice. The choice would be the closest of the three to the estimated one and most likely would give comparable results. This conclusion is based on the empirical evaluations in Tables 3 and 4, where the best of the three is only slightly less efficient than  $\hat{Y}_*$  coupled with the fact that an estimated predictor is never as efficient as  $\hat{Y}_*$ .

JENKINS (1934) did not find a great deal of difference among the correlations of the three single cross predictors with double cross performance. He was estimating the following correlations,

$$\rho_1 \tilde{Y} = \rho_{1Y} \rho_{\tilde{Y}Y}, \quad \rho_2 \tilde{Y} = \rho_{2Y} \rho_{\tilde{Y}Y}, \quad \rho_3 \tilde{Y} = \rho_{3Y} \rho_{\tilde{Y}Y}. \quad (55)$$

One can see from the tabulations in Tables 3 and 4 that there are many situations for which the three are not very different, and not enough different to be distinguished from data. Some paradoxes are brought out in the Tables. The variance of an estimated  $\rho_i \tilde{Y}$  may be reduced by increasing  $p$  or increasing replications represented by increasing  $\beta$  in the Tables. However, there are genetic situations for which the three  $\rho_i \tilde{Y}$ 's become more alike with an increase in  $p$  and/or  $\beta$ . Whether the increased similarity is more than offset by the decrease in variance of the estimates requires further exploration, but there is the suggestion of not being offset. Certainly, reasonably good data are required for distinguishing among the predictors. The double cross values should be well enough estimated that if another similar but environmentally independent set of estimates,  $\tilde{Y}'$ , were available, the correlation between the two,

$$\rho_{\tilde{Y}\tilde{Y}'} = \rho_{\tilde{Y}Y}^2, \quad (56)$$

would be reasonably high.

EBERHART (1964) and EBERHART et al. (1964) considered the use of three-way crosses in addition to single crosses in the prediction of double crosses. From a fixed sample viewpoint, they were principally concerned with minimum genetic biased predictors. The prediction model (4) can be augmented to include three-way cross information. Taking into account the different ways in which three-way crosses relate to double crosses will add at least six new variables,  $Z$ 's, to the model (4). These additional variables can probably be formulated into an uncorrelated set as was done for the single crosses, although some will be correlated with the  $Z$ 's pertaining to the single crosses. Minimization of the prediction error variance and maximization of the gain is accomplished in the same manner as was done for only single cross prediction. The inclusion of three-way cross information properly weighted will always improve the efficiency of prediction over the use of single crosses alone. The amount of improvement can be evaluated empirically for various parametric situations.

Only a single replicated experiment was considered for simplicity, but only slight modifications are needed to include environmental factors and genetic by environmental interactions. For an experiment at  $l$  environments for example, one replaces  $M_3/r$  and  $M_2/r$  in (6) with

$$\left. \begin{aligned} \frac{M_3}{rl} &= \frac{\sigma^2}{rl} + \frac{\sigma_{gE}^2}{l} + (p-2) \sigma_g^2, \\ \frac{M_2}{rl} &= \frac{\sigma^2}{rl} + \frac{\sigma_{sE}^2}{l} + \sigma_s^2, \end{aligned} \right\} \quad (57)$$

where  $\sigma_{gE}^2$  and  $\sigma_{sE}^2$  are components of variance for general and specific by environments. The gene effects and variances and covariances of relatives are always defined to be free of environments and of interactions of genes with environments, and any estimates of them should take this into account. A corresponding modification must be made in  $\sigma_{\tilde{Y}}^2$  for the inclusion of different environments. An extension to multiple environmental classifications is straightforward, just more complicated.



Prediction and gain have been considered from an internal sample standpoint. By this is meant that only contrasts among the  $Y$ 's, and as predictors only contrasts among the  $X$ 's, have been considered. In such case the means, both of double and single crosses, cancel out whether from a random or fixed sample viewpoint. The question arises as to how to develop comparative predictors for separate samples of single crosses in different experiments. The means of the samples will differ both genetically and environmentally. In most cases the experiments are in different years and often in different places. Consequently, environmental differences are much more than those reflected by averages of internal environmental effects. One solution is to augment the prediction equation (4) with  $b_4 Z_4$  where  $Z_4$  is a deviation of an experimental mean of single crosses from the mean of all experiments. The relative weight,  $b_4$ , depends mainly on the proportions of genetic and environmental variances among experiments for which little information is generally available. If the environmental differences are relatively large then  $b_4$  will be relatively small, and the comparisons are very nearly among internal predictions. Internal predictions are probably the best and safest solution in practice.

Having decided on a particular predictor,

$$\hat{Y} = b_1 (X_n - X_{\bar{p}}) + b_2 (X_{\bar{p}} - 2 X_l + X_0) + b_3 (X_l - \bar{X}), \quad (58)$$

the equation may be rearranged algebraically to simplify calculations. One alternative which is computationally simple is

$$\hat{Y} = b'_1 X_p + b'_2 X_{\bar{p}} + b'_3 X_l + b'_4 \bar{X}, \quad (59)$$

where

$$\left. \begin{aligned} b'_1 &= \frac{-b_1}{2}, & b'_2 &= \frac{b_1}{2} + b_2 \frac{(p-1)(p-2)}{(p-4)(p-5)}, \\ b'_3 &= b_3 - b_2 \frac{2(p-1)^2}{(p-4)(p-5)}, \\ b'_4 &= b_2 \frac{p(p-1)}{(p-4)(p-5)} - b_3. \end{aligned} \right\} \quad (60)$$

The last term,  $b'_4 \bar{X}$ , is a constant and may be dropped in ranking the double crosses, and also one of the coefficients may be standardized to one. For example,

$$b''_1 X_p + X_{\bar{p}} + b''_3 X_l, \quad (61)$$

where

$$b''_1 = \frac{b'_1}{b'_2}, \quad b''_3 = \frac{b'_3}{b'_2}, \quad (62)$$

will rank the double crosses the same as (58).

## Appendix

### Variances among Double Crosses

From (38) the experimental value for double cross  $AB \cdot CD$ , for example, is

$$\tilde{Y}_{AB \cdot CD} = Y_{AB \cdot CD} + \bar{\epsilon}_{AB \cdot CD}, \quad (63)$$

and the experimental error is the mean of the plot errors for the double cross,

$$\bar{\epsilon}_{AB \cdot CD} = \frac{\sum \epsilon_{AB \cdot CD}}{r} \quad (64)$$

with variance

$$\sigma_{\bar{\epsilon}}^2 = \frac{\sigma_{\epsilon}^2}{r}. \quad (65)$$

We shall use the following measure of variance,

$$\sigma_{\tilde{Y}}^2 = E(\tilde{Y}_{AB \cdot CD} - \bar{\tilde{Y}})^2, \quad (66)$$

where  $\bar{\tilde{Y}}$  is the mean of all double crosses under consideration. This measure of variance brings into consideration the number of double crosses and the correlations among them in keeping with the treatment of the single cross predictors. The variance reduces to that among the true genotypic values of the double crosses by making  $\sigma_{\epsilon}^2 = 0$ .

Expanding (66),

$$E(\tilde{Y}_{AB \cdot CD} - \bar{\tilde{Y}})^2 = E(\tilde{Y}_{AB \cdot CD}^2) - 2E(\tilde{Y}_{AB \cdot CD} \bar{\tilde{Y}}) + E(\bar{\tilde{Y}}^2), \quad (67)$$

it can be shown that

$$E(\tilde{Y}_{AB \cdot CD} \bar{\tilde{Y}}) = E(\bar{\tilde{Y}}^2), \quad (68)$$

so that

$$\sigma_{\tilde{Y}}^2 = E(\tilde{Y}_{AB \cdot CD}^2) - E(\bar{\tilde{Y}}^2). \quad (69)$$

Evaluating the first expectation,

$$E(\tilde{Y}_{AB \cdot CD}^2) = \mu^2 + C_{AB \cdot CD, AB \cdot CD} + \sigma_{\epsilon}^2, \quad (70)$$

where  $C_{AB \cdot CD, AB \cdot CD}$  is the covariance of the genotypic values of the double crosses with themselves, which is the variance among the true values of unrelated double crosses, but will be left in covariance terminology in accord with the remaining expectations. Turning now to the second expectation in (69), taking into account all possible types of related hybrids in the mean of the  $p(p-1)(p-2) \times (p-3)/8$  double crosses, and assuming of course that all experimental errors are uncorrelated in any way, the following result is obtained

$$\begin{aligned} E(\tilde{Y}_{AB \cdot CD} \bar{\tilde{Y}}) &= \mu^2 + \left[ \sigma_{\epsilon}^2 + 1 \right. & C_{AB \cdot CD, AB \cdot CD} \\ &+ 2 & C_{AB \cdot CD, AC \cdot BD} \\ &+ 4(p-4) & C_{AB \cdot C-, AB \cdot C-} \\ &+ 8(p-4) & C_{AB \cdot C-, AC \cdot B-} \\ &+ 4(p-4)(p-5) & C_{A \cdot -B-, A \cdot -B-} \\ &+ 4(p-4)(p-5) & C_{A \cdot -B-, AB \cdot -} \\ &+ (p-4)(p-5) & C_{AB \cdot --, AB \cdot -} \\ &+ 2(p-4)(p-5)(p-6) & C_{A \cdot --, A \cdot -} \\ &+ \frac{(p-4)(p-5)(p-6)(p-7)}{8} & C_{-----} \left. \right] / \\ &\frac{p(p-1)(p-2)(p-3)}{8}. \end{aligned} \quad (71)$$

The designations of the covariances show the number of parental lines common and the manner in which they are common. The coefficients of the covariances add to the total number of double crosses, the last covariance with no lines common being included for completeness and is always zero. More details are given by COCKERHAM (1961) and RAWLINGS and COCKERHAM (1962). With the assumptions generally required in expressing the covariances of relatives

in terms of genetic variances (16) the additive and dominance coefficients are as follows:

$$\left. \begin{array}{l} C_{AB \cdot CD, AB \cdot CD} \\ C_{AB \cdot CD, AC \cdot BD} \\ C_{AB \cdot C-, AB \cdot C-} \\ C_{AB \cdot C-, AC \cdot B-} \\ C_{A \cdot \cdot B-, A \cdot \cdot B-} \\ C_{A \cdot \cdot B-, AB \cdot \cdot \cdot} \\ C_{AB \cdot \cdot \cdot, AB \cdot \cdot \cdot} \\ C_{A \cdot \cdot \cdot \cdot, A \cdot \cdot \cdot \cdot} \\ C_{\cdot \cdot \cdot \cdot \cdot \cdot, \cdot \cdot \cdot \cdot \cdot \cdot} \end{array} \right\} \begin{array}{l} \frac{\alpha}{2} \\ \frac{1}{2} \\ \frac{3}{8} \\ \frac{3}{8} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{8} \\ 0 \end{array} \left. \begin{array}{l} \frac{\delta}{4} \\ \frac{1}{8} \\ \frac{1}{8} \\ \frac{1}{16} \\ \frac{1}{16} \\ 0 \\ 0 \\ 0 \\ 0 \end{array} \right\} \quad (72)$$

Putting (70) and (71) together as in (69) and substituting the coefficients (72) into the covariances as in (16) leads to the appropriate expressions for  $\sigma_Y^2$  and  $\sigma_e^2$ ,

$$\left. \begin{array}{l} \sigma_Y^2 = \left(1 - \frac{8}{p(p-1)(p-2)(p-3)}\right) \frac{\sigma_e^2}{r} + \sigma_Y^2, \\ \sigma_Y^2 = \left(1 - \frac{4}{p}\right) \frac{\sigma_A^2}{2} + \left(1 - \frac{8}{p(p-1)}\right) \frac{\sigma_D^2}{4} \\ + \left(1 - \frac{p+8}{p(p-1)}\right) \frac{\sigma_{AA}^2}{4} + \dots \end{array} \right\} \quad (73)$$

The error variance,  $\sigma_e^2$ , is not exactly comparable to that,  $\sigma^2$ , for the single crosses. The error variance for single crosses, in the absence of competitive effects, can be decomposed into

$$\sigma^2 = \frac{\sigma_1^2}{n} + \sigma_2^2, \quad (74)$$

where  $\sigma_1^2$  is the environmental variance among  $n$  plants in a plot and  $\sigma_2^2$  is the plot component of environmental variance. In contrast,

$$\sigma_e^2 = \frac{\sigma_1^2 + C_{S2} - C_{AB \cdot CD, ABCD}}{n} + \sigma_2^2, \quad (75)$$

where  $C_{S2} - C_{AB \cdot CD, ABCD}$  represents the genotypic variance among individuals of the same double cross

in a plot, and  $C_{S2}$ , given in (17), represents the total genotypic variance.

### Zusammenfassung

Für die Vorausschätzung von Doppelkreuzungsbastarden aus Einzelkreuzungshybriden wird eine einheitliche Theorie entwickelt, die sowohl genetische als auch experimentelle Bedingungen berücksichtigt. Die Methode ist der für die Berechnung von Selektionsindizes analog. Es wird die Beziehung des Vorhersagemodells zum genetischen Modell erläutert. JENKINS' (1934) drei Einzelkreuzungs-Schätzwerte, der beste Einzelkreuzungs-Schätzwert und die Selektion auf der Grundlage der Doppelkreuzungs-Schätzungen werden für ein additives und Dominanz-Modell mit variierenden Verhältnissen der experimentellen Fehlervarianz und unterschiedlicher Anzahl von Bastarden empirisch miteinander verglichen. Der Unterschied zwischen fixierter und zufälliger genetischer Stichprobenmethode wird hinsichtlich der Vorhersage besprochen.

### References

1. COCKERHAM, C. CLARK: Implications of genetic variances in a hybrid breeding program. *Crop Science* 1, 47-52 (1961).
2. COCKERHAM, C. CLARK: Estimation of genetic variances. *Statistical Genetics and Plant Breeding*. National Academy of Sciences-National Research Council Publ. 982, 53-94 (1963).
3. COCKERHAM, C. CLARK: Group inbreeding and coancestry. (Submitted to *Genetics*) (1967).
4. EBERHART, S. A.: Theoretical relations among single, three-way, and double cross hybrids. *Biometrics* 20, 522-539 (1964).
5. EBERHART, S. A., W. A. RUSSELL, and L. H. PENNY: Double cross hybrid prediction in maize when epistasis is present. *Crop Science* 4, 363-366 (1964).
6. HENDERSON, C. R.: Selection index and expected genetic advance. *Statistical Genetics and Plant Breeding*. National Academy of Sciences-National Research Council Publ. 982, 141-163 (1963).
7. JENKINS, MERLE T.: Methods of estimating the performance of double crosses in corn. *J. Amer. Soc. Agronomy* 26, 199-204 (1934).
8. LUSH, JAY L.: Family merit and individual merit as bases for selection. *Amer. Naturalist* 81, 241-261; 362-379 (1947).
9. PATEL, R. M.: Selection among factorially classified variables. Ph. D. Thesis, North Carolina State University at Raleigh (1962).
10. RAWLINGS, J. O., and C. CLARK COCKERHAM: Analysis of double cross hybrid populations. *Biometrics* 18, 229-244 (1962).
11. SPRAGUE, G. F.: Chapter V: Corn breeding. In: *Corn and Corn Improvement*. New York, N.Y.: Academic Press Inc. 1955.
12. WILLIAMS, J. S.: Some statistical properties of a genetic selection index. *Biometrika* 49, 325-337 (1962).

## A Comprehensive Breeding System\*

S. A. EBERHART, M. N. HARRISON and F. OGADA<sup>1</sup>

Maize Research Section, Kitale, Kenya

**Summary.** An outline of a comprehensive breeding system developed and used by the Kenya Maize Research Section is presented. This system has four main phases:

1. Evaluation of local and exotic varieties so that the breeding material is the best available.
2. Compositing the selected breeding material into two or more populations or varieties in such a manner that

each population has considerable genetic variation for the traits requiring improvement and that the crosses of these populations will show heterosis.

3. Recurrent selection in each population to increase the frequency of favorable genes so that the populations and population crosses are improved with each cycle of selection.

4. Release of a commercial variety of one of the following forms: (a) the cross of two populations as a variety cross hybrid; (b) single, three-way or double cross hybrids from inbred lines developed from the elite material after each cycle of selection; or (c) a synthetic variety derived from the advanced generation of the population cross in areas where hybrid production is not yet feasible.

Preliminary results are presented to indicate the improvement possible in maize by use of this system. Its possible extension to other crops is also briefly discussed.

\* Dedicated to Dr. GEORGE F. SPRAGUE on the occasion of his 65<sup>th</sup> birthday.

<sup>1</sup> Research Geneticist, Agricultural Research Service, U.S. Department of Agriculture/U.S. Agency for International Development/East African Agriculture & Forestry Research Organisation; Senior Maize Research Officer and Maize Breeder, Maize Research Section, Kenya Ministry of Agriculture.