HILGARDIA

A Journal of Agricultural Science Published by the California Agricultural Experiment Station

VOLUME 16

FEBRUARY, 1944

NUMBER 1

CONTENTS

BIASES ENCOUNTERED IN LARGE-SCALE YIELD TESTS

O. C. RIDDLE and G. A. BAKER

UNIVERSITY OF CALIFORNIA · BERKELEY, CALIFORNIA

HILGARDIA

A Journal of Agricultural Science Published by the California Agricultural Experiment Station

VOL. 16

FEBRUARY, 1944

No. 1

BIASES ENCOUNTERED IN LARGE-SCALE YIELD TESTS^{1,3}

O. C. RIDDLE⁸ and G. A. BAKER⁴

INTRODUCTION

CERTAIN DIFFICULTIES in interpreting the results of the ordinary analysis of variance when applied to yield tests of genetically similar wheat strains, derived through backcrossing, prompted a critical examination of the data. These data indicated a bias inherent in any large-scale experiment that imposes an inflexible design upon a soil whose fertility may fluctuate markedly within short distances. Because of this bias, the usual analysis of variance tends to overestimate significance grossly when the number of varieties is large and the productivity levels of the soil change rapidly and erratically (as at Davis). One may briefly describe the bias by saying that the inflexible design of the experiment tends to subtract too little from the naturally high-yielding plots and too much from the naturally low-yielding plots in attempting to correct for differences in soil productivity. This has a spreading effect on the part of the variation that is labeled "varietal differences" and thus causes a serious overestimation of significance.

In a conventional analysis of variance, the natural variation of an experiment is arbitrarily partitioned into categories according to a preconceived mathematical model. These categories are labeled "variation due to varieties," "variation due to soil productivity," and so on. If the experiment is in exact accord with the model, the labels are accurate. If the experiment is not as called for by the model, the labels are misleading: for instance, the category labeled "variation due to varieties" may contain some of the variation due to soil productivity.

A mathematical model may be pleasing and beautiful to its creator or users. If, then, nature does not conform to the model, workers having limited firsthand experience with the vagaries of biological material may even feel that nature has erred and should be corrected. Baten, Northam, and Yeager (2)⁵,

¹ Received for publication May 14, 1943.

² Results of a cooperative study conducted by the United States Bureau of Plant Industry Division of Cereal Crops and Diseases, and the University of California Division of Agronomy and Division of Mathematics and Physics, Davis, California.

Instructor in Agronomy and Junior Agronomist in the Experiment Station. Formerly Agent, United States Department of Agriculture.

Assistant Professor of Mathematics and Assistant Statistician in the Experiment Station.

for example, in a recent issue of *Journal of the American Society of Agronomy*, throw out the observed yields on two strains of tomatoes because those yields are not in accord with their mathematical model. The observed yields are replaced by computed yields based on the other observed yields. There is another strong temptation: a model that has proved useful in a restricted realm may be unquestioningly extrapolated to new and unexplored situations. This has been done by workers in all branches of science. One main purpose of this paper is to point out the danger of such extrapolation.

To construct an adequate mathematical model of some portion of nature, one must have extensive and accurate data on all situations to which the model is to apply. Such data are not now available. This paper aims to stimulate the collection of data that will provide for an improved model for yield trials. It may help to show that such data are necessary.

GENETIC RELATIONSHIP OF MATERIAL

As pointed out by Suneson and his co-workers (17), 182 F_s strains of the pedigree (Martin × White Federation⁶) × (Hope × White Federation⁵) and 157 F_s strains of (Martin × Baart⁷) × (Hope × Baart⁵) were bulked to produce White Federation 38 and Baart 38, respectively. These two new wheat varieties have been shown (17) to be essentially like their prototypes except for the incorporated resistance to bunt, or stinking smut (*Tilletia tritici*), and to stem rust (*Puccinia graminis* var. *tritici*). From the strains mentioned above, selections were made at random for the yield tests reported in this paper.

An understanding of the genetic relationship of these strains is important in consideration of the results of this test and in the extension of the implications therefrom to other methods of testing similar or unrelated material.

As Briggs (3, 4) has pointed out, "the proportion of homozygous individuals in any backcross generation is the same as would result from an equal number of selfed generations." This proportion may be calculated from the equation

per cent homozygosity =
$$\left(\frac{2^m - 1}{2^m}\right)^n \times 100,$$
 (1)

where m is the number of generations of backcrossing and n is the number of heterozygous factor pairs in the original cross.

As Jones (8, p. 331) explains, "the proportion of complete homozygotes to the different classes of heterozygotes in any generation" can be obtained by expanding the binomial

$$[1 + (2m - 1)]n, (2)$$

where *m* and *n* are as above.

If we assume that the wheat parents used to develop the strains under test differ by 21 factor pairs (that is, a 1-factor difference on each of their chromosome pairs) and that the genotypes have been randomly sampled in the course of backcrossing, we can approximate the degree of homozygosity of the wheat strains in question. On the basis of these two assumptions, we can show by means of equations 1 and 2 that in the F_3 of the sixth backcross, 82.3 per cent of the plants are homozygous for the recurrent parent genotype, and that an additional 13.5 per cent differ by only 1 factor.

Admittedly, the parents differ by more than 21 factor pairs; but calculations accounting for all factor differences and the effects of linkage are not possible. Certain considerations suggest, however, that the estimates of homozygosity given above may approach the true situation. As the fraction $\frac{2^m-1}{2^m}$

(from equation 1) approaches unity, n (the number of heterozygous factor pairs) will change that value but little. Through the mechanism of crossing over and random assortment of chromosomes, repeated backcrossing facilitates recovery of the complete chromosome complement of the recurrent parent except for a small segment from the nonrecurrent parent on which the gene being selected for is located. In developing these strains, rigid selection toward the recurrent parent phenotype was practiced in early backcross generations; this hastens return to the recurrent parental genotype and reduces the unfavorable effects of linkage.

Actually, under the conditions of these tests the strains were morphologically indistinguishable, except strain 1441, which averaged 2 inches taller than the others. No other differences in growth habit were observed at any stage, and there was no detectable difference in reaction to any disease among the strains, although the original Baart and White Federation grown with them were attacked by rust.

Although the strains are not genetically identical, the degree of similarity is clearly such as to require a critical test of significance to detect any possible differences in yielding ability.

METHODS AND DESIGN

Twenty-nine strains of Baart 38 and thirty-four strains of White Federation 38 were chosen at random for yield testing. These, together with the disease-susceptible prototypes, were set up in separate experiments in each of the years 1939 and 1940. The design employed for these tests was a modified Latin square suggested by "Student" (16) and by Snedecor (13, p. 38) and used extensively by Pope (11) and others (15, 21). There were five replications divided into five columns superimposed upon and situated at right angles to the replications. The strains were grown in single 16-foot rows. They were completely randomized except for the double restriction that each strain occurred once (and only once) in each replication and each column. The same randomization but totally different fields were used for the tests in the two years. The high and low extremes from the 1939 and 1940 tests of both the Baart 38 and the White Federation 38 strains, each with its susceptible prototype, were tested in separate 6×6 Latin squares in 1941 using single 16-footrow plots.

The data were analyzed by use of the analysis of variance, testing for significance by F (14, p. 184). Least significant differences above and below the general mean were established by the use of t (14, p. 58) for the appropriate degrees of freedom. The values of $\frac{\text{range}}{S.E.}$ expected from random sampling in a

normal homogeneous population of sample size N were obtained from Snedecor (14, p, 89).

Yates (19) has criticized the modified Latin square as being subject to a biased estimate of error. Repeated reference to this criticism in the literature (5, 20) condemns the design in favor of the more complex incomplete block

designs. The present writers do not defend the modified Latin square design nor advocate its use; but they feel that information can be gained from the data in these experiments in which the modified Latin square design was used.

There is evidence that the error estimates in these experiments are not biased in the sense of Yates's (19) criticism. The F values (table 4) due to strain differences are greater than required for the 1 per cent level of significance in all cases for the 1939 and 1940 modified Latin square tests. If these F values are large only because of a biased estimate of error, then the estimates of error variance must be considered *too small in all cases*. We may compare certain variances as a test of Yates's bias. The degrees of freedom for the Baart 38 tests, for instance, may be set out from what Fisher (7) calls a topographical standpoint as:

	replications	4
Between plots	columns	4
	replications \times columns	16
Within plots	ζ -	125
Total		149

If the interaction variance of replications \times columns is the same as the withinplot variance, then no Yates's bias can exist. A comparison of these variances for the 1939 and 1940 experiments in table 1 shows no significant nor consistent differences when tested by F and therefore no indication of Yates's bias.

1 1.39	1.72
5 1.05	2.04
6 1.02	2.04
4 1.38	1.71
1	6 1.02 4 1.38

 TABLE 1

 Comparison of Variances Indicating No Yates's Bias in the Data

Yates's bias, if it exists, seems to be sometimes in one direction and sometimes in the other, hence resembles the biases considered by Welch (18). Such a bias can be expected to balance out in the long run; one could allow for it by slightly changing the probability levels at which significance is accepted or rejected.

EXPERIMENTAL RESULTS

Yield data for the Baart 38 and White Federation 38 component strains tested in 1939, 1940, and 1941 are given in tables 2 and 3, respectively. The several strains are listed in descending order of average yields for 1939 and 1940 combined. Row numbers per replication are listed for each strain in recording the 1939 and 1940 results, to facilitate additional treatment of the data if desired. Table 4 summarizes the analysis of variance for each experiment, and gives pertinent statistical constants.

Indicated Significances.—In all 1939 and 1940 tests, F values due to differences in mean yields of strains exceed those required at the 1 per cent level of significance (see table 4). Standard errors were calculated; and minimum

Yields and Row Numbers for 1939–1940 Modified Latin Square Yield Trials, and Mean Yields for 1941 Latin Square Test of Baart and of Randomly Selected Component Strains of Baart 38 TABLE 2

Strain mean 1941		401.0 419.5 365.7 365.7	377.2‡
1939–1940 com- bined mean		606 5794 (200 5714 (200 57	407.11
Strain mean 1940		630.0 645.3 645.3 645.3 645.3 645.3 645.3 655.6 645.3 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 655.6 6 655.6 6 6 6	459.21
Strain mean 1939		582* 551 551 551 555 555 555 555 555 555 55	64 4
Λ	1940 yield	630 612 613 613 613 613 653 573 573 573 573 573 573 573 573 573 5	480
plicatio	1939 yield	560 5550 5550 5550 5550 5550 5550 5550	010
1940 Rej	Row no.	&\$\$\$ \$ \$	4 4 :
n: 1939- 1 IV	1940 yield	628 537 537 537 537 537 537 537 537 537 537	357
plication	1939 yield	445 450 3330 3345 3345 3345 3345 3355 3375 544 450 3375 544 450 3375 544 450 3375 544 450 3375 544 450 3375 544 5450 3375 5450 3375 5450 3375 5450 5450 5455 5455 5455 5455 5455 54	445
. per rel	Row no.	8888898988898469449469889489889889898989	1 23 :
d row number lication III	1940 yield	782 755 555 555 555 555 555 555 555 555 55	504
	1939 yield	790 7455 7455 7455 7455 7455 7455 7455 745	930 · · ·
t row al Rep	Row no.	88748874888888664448848855558888486488	<u>8</u> :
r 16-foo n II	1940 yield	477 5390 5390 5390 5390 5391 5391 5393 5593 5593 5593 5593 5593	430
ams pe	1939 yield	540 5440 5440 6415 5550 5550 5550 5550 5550 5550 5550 5	066
ld in gr Re	Row no.	&#&#&&&&&&&&&&&&&&&&&&&&&&&&&&&&&&&&&&</td><td>fc :</td></tr><tr><td>Yie n I</td><td>1940 yield</td><td>525133 52515 52515 52515 52515 52515 52515 52515 52515 52515 5</td><td>4Z0 </td></tr><tr><td>plicatio</td><td>1939 yield</td><td>575 540 550 5540 5540 5540 5555 5455 5555 5555 5555 5555 5555 5555 5555 5555</td><td>435</td></tr><tr><td>Re</td><td>Row no.</td><td>444997888888888888888888888888888888888</td><td>49</td></tr><tr><td>Strain no.</td><td></td><td>5637 2755 2756 4673 4673 4673 4613 4614 4630 33541 33541 33541 33541 33541 5485 5485 5485 5485 5585 5585 5585 5</td><td>3769 . Baart 38‡</td></tr></tbody></table>	

Significantly higher than the general mean. The general mean was 512.7 in 1939, 550.7 in 1940, and 531.7 for the combined yields.
 Significantly lower than the general mean.
 The mixture, Baart 38, was not grown in 1939 and 1940 since the object was to determine if Baart 38 could be separated into parts that yield differently.

Yields and Row Numbers for 1939-1940 Modified Latin Square Yield Trials, and Mean Yields for 1941 Latin Square Test of White Federation and of Randomly Selected Component Strains of White Federation 38 TABLE 3

Strain mean 1941			997.7 397.7 3894.8 394.8 379.2 379.2 405.77
1939–1940 com- bined mean		mean	636 9 636 9 618 9 61
	Strain mean 1940		672 672 672 672 672 672 672 672
	Strain mean 1939		000 014 014 014 014 014 019 019 019 019 019 019 019 019
	n V	1940 yield	665 625 625 625 625 625 647 647 647 652 6538 6538 6518 6518 6518 6518 6518 6528 6528 6528 6528 6528 6528 6528 652
	lication	1939 yield	640 550 550 550 550 550 550 550 550 550 5
1940	Rej	Row no.	8022222224800500884822222883995389538953895389538595555555555
ı: 1939-	١V	1940 yield	629 629 640 655 555 555 555 555 557 657 657 657 657
olication	lication	1939 yield	560 5540 5540 5540 5540 5550 5575 5575 557
per re	Rep	Row no.	: 233523252352352323232322449028611932252
number	III	1940 yield	710 710 723 598 597 597 597 507 501 501 503 557 557 557 557 557 557 557 557 557 55
lication	lication	1939 yield	650 650 710 633 6335 6335 6335 6335 6335 6335 6335
t row ai	Rep	Row no.	: 888242 2699 - 2011 238 200 238 238 28 28 28 28 28 28 28 28 28 28 28 28 28
r 16-foo	n II	1940 yield	745 745 745 745 745 745 745 745
ams pe	plicatio	1939 yield	665 665 7040 610 610 610 610 605 605 605 605 605 605 605 605 605 60
ld in gr	Re	Row no.	: \$\$\$52258\$\$ 33354 558\$\$\$5233555555555555555555555555555555
Yie	n I	1940 yield	 3555208 55209 <l< td=""></l<>
	plicatic	1939 yield	$\begin{array}{c} 488\\ 488\\ 550\\ 550\\ 550\\ 550\\ 550\\ 550\\ 550\\ 5$
Re	Re	Row no.	8488282882882882888288828882882882882882
Strain no.			2029 153 153 153 243 244 245 249 245 249 249 167 257 167 167 167 167 167 167 167 16

Significantly higher than the general mean. The general mean was 584.2 in 1939, 577.1 in 1940, and 580.6 for the combined yields. Fignificantly lower than the general mean. † The mixture, White Federation 38, was not grown in 1339 and 1940 since the object was to determine if White Federation 38 could be separated into parts that yield differently.

TABLE 4

SUMMARY OF VARIANCE ANALYSIS AND STATISTICS OF BAART 38 AND WHITE FEDERATION 38 STRAINS TESTED FOR YIELD, 1939–1941

Experiment and source	Degrees	Sum of	Mean	F value			
of variation	freedom	squares	square	Actual	5 Per cent	1 Per cent	
a. Baart 38 strains—1939: Between mouns of replicator	4	594 649	121 169				
Between means of replicates	4	024,040	131,102				
Between means of columns	- 4	307,231	91,813	0.10	1	1 00	
From	119	200,139	0,901	2.10	1.0/	1.89	
Enfor		308,291	3,288				
Total	149	1,460,329					
b. Baart 38 strains-1940:							
Between means of replicates	4	202,727	50,682	1			
Between means of columns	4	40,766	10,192				
Between means of strains	29	283,937	9,791	2.02	1.57	1.89	
Error	112	544,176	4,859				
Total	149	1,071,606					
c. Baart 38 strains-1941:							
Between means of replicates	5	25,432	5,086				
Between means of columns	5	16,783	3,357				
Between means of strains	5	18,700	3,740	0.71	2.71	4.10	
Error	20	105,850	5,292				
Total	35	166,765					
d. White Federation 38 strains-1939:							
Between means of replicates	4	894,616	223,654				
Between means of columns	4	246.805	61,701				
Between means of strains	34	186,975	5,499	2.15	1.55	1.85	
Error	132	337,509	2,557				
Total	174	1,665,905					
e. White Federation 38 strains-1940:							
Between means of replicates	4	75,467	18,867				
Between means of columns	4	9,565	2,391				
Between means of strains	34	438,644	12,901	3.46	1.55	1.85	
Error	132	491,588	3,724				
Total	174	1,015,264					
f. White Federation 38 strains-1941:							
Between means of replicates	5	80,757	16,151				
Between means of columns	5	4,422	884				
Between means of strains	5	15,430	3,086	1.02	2.71	4.10	
Error	20	60,613	3,031				
Total	35	161,222					

Statistics		Experiment						
	a	b	c	d	е	f		
1. General mean (gms. per 16 ft. row)	512.7	550.7	397.5	584.2	577.1	399.2		
2. S.E. of a single plot (gms.)	57.3	69.7		50.6	61.0			
3. Least significant difference at 5 per cent level between								
any strain mean and general mean (gms.)	51.4	62.6		45.2	54.6			
4. Range in strain means $/S.E.$ of a strain mean:								
Observed	5.7	6.7		7.8	10.5			
Expected (approximate)	4.1	4.1		4.1	4.1			
5. Coefficient of variability (per cent)	5.0	5.7		3.9	4.7			
				1 1		1		

significant differences, at the 5 per cent level, above and below the general mean were established. The general mean yield of all strains was used as the reference point for judging significance because, if significant differences did exist, the strains higher in yield than the general mean of all strains would be of greatest agronomic value.

As indicated in tables 2 and 3, certain strains were found to differ significantly from the general mean in both 1939 and 1940 and in the combined 1939 and 1940 results. In addition, two criteria suggest that the observed mean differences are not of the order expected from random sampling in a homogeneous population : first, the high values of F; second, the high values of range in mean yields divided by standard error of the mean (table 4).

Difficulties of Interpretation.—Considering the genetic relationship of the strains tested, we might expect that their mean yields would not be significantly different. Yet, as mentioned above, significant differences above and below the general mean are indicated, and they are such as to deserve consideration agronomically. Furthermore, if differences do exist, the strains might be expected to maintain somewhat the same relative positions in repeated tests. Certainly such characters as growth habit at any stage, plant height, or date of maturity showed no observable difference that would suggest a differential response to environment. Moreover, of the strains indicated as significantly different from the general mean in tables 2 and 3, only one strain, 5037 (a Baart 38 strain), holds the same position for the two years, 1939 and 1940, in having a yield significantly greater than the general mean. Two strains, namely 5129 (a Baart 38 strain) and 1617 (a White Federation 38 strain), reversed their position from significantly lower in 1939 to significantly higher than the general mean in 1940. All other strains that were significantly different in one or the other of the two years were not significantly different in the alternate year. As will be observed in tables 2 and 3, the Baart and White Federation prototypes did not differ significantly from the general mean in 1940 even though they were attacked by stem rust.

So far in the analysis, two reactions have been indicated that are not in accord with expectation if the genetic relationship of the strains agrees with that expressed under the section on "Genetic Relationship of Material." These reactions are, namely, significant differences in yield, and failure to exhibit comparable relative yield responses in two different years.

The correlation of one year's yields with another year's depends upon the variance of the true yields and upon the experimental-error variance. If for the two years the variances of the true yields are the same, if the strains maintain their relative positions, and if the error variance is the same for every plot, then the expected correlation between the yields for the two years is a value equal to the variance of the true yields divided by the quantity true-yield variance plus strain-mean variance; hence it is less than 1. Conceivably, the spread among the true yields might be large enough to insure detection four times out of four, and yet small enough relative to the error of the experiment so that the correlation between yields for two years might be small. Neyman's tables⁶ for the probability of failing to detect differences between

^e Supplied in manuscript form by Professor J. Neyman of the Statistical Laboratory, University of California, Berkeley. Feb., 1944]

yields when they actually exist show that this situation is not possible for these data.

The 1941 6×6 Latin square yield test, sampling the 1939 and 1940 high-low extremes, constituted a more critical attempt to determine whether the strains actually differ in yielding ability. If true yield differences exist, they should certainly be represented in those extremes indicated by previous tests to be significantly different. As shown in the analysis of variance applied to the 1941 data (table 4), the *F* values due to differences in mean yields of previously indicated high-low strains are not significant. The results of this test do not indicate differences in yield of the strains tested. Neyman's tables indicate that this test was sufficient to detect, with a probability greater than 0.8, differences of the order necessary to account for the observed *F* values if the usual mathematical model is assumed.

Considering these difficulties, we may well examine critically the statistical methods used and the assumptions on which those statistics are based.

BASIC ASSUMPTIONS IN ANALYSIS OF VARIANCE

The derivation of the analysis-of-variance technique is based on four assumptions: (1) that the productivity levels of the plots assigned to any variety are independent of those assigned to any other; (2) that the estimates of individual plot yields are normally distributed about the "true" plot yield; (3) that the distribution of yield estimates for every plot has the same variance; (4) that in yield trials the productivity levels follow some prescribed law.

The work of Neyman (10) and McCarthy (9) relates to the first assumption. According to Neyman, the Latin square design may often indicate significance between hypothetical "varieties" in uniformity trials, partly because of unequal correlations between fertility levels of the plots assigned to the different "varieties." McCarthy shows further that these unequal correlations, when the varieties are tested in randomized blocks, may cause a serious overestimation of significance—that is, may indicate significance where none exists.

According to the second and third assumptions, the estimates of individual plot yields are normally distributed about the "true" plot yield with equal variances. In this connection the work of Baker (1) and Salmon (12) may be cited. Baker shows that the distribution of the estimates of a "true" plot yield is usually skewed one way or the other and that the variance of the distribution of estimates depends on the variations of the fertility within the plot. Sometimes the nonnormality of the distribution may be such as to cause a serious overestimation of significance. Baker shows further that adjacent plots of 15 square feet cannot be assumed to have the same fertility levels. Salmon and many others have discussed unequal variance as affecting the results of the analysis of variance.

The importance of assumption 4 is well recognized by some authorities (see Neyman [10]), but has been generally overlooked or deëmphasized by many workers concerned with yield trials. The present data show clearly the importance of the failure of conventional designs to prescribe a sufficiently flexible law for fertility levels.

Failure of the data to comply with any one of the assumptions on which the statistic is based may result in misleading or invalid conclusions.

CRITICAL EXAMINATION OF DATA

The nature and extent of the material tested and the relative simplicity of the design employed in these tests facilitate a critical examination of the data from the standpoint of the validity of the four assumptions mentioned above.

Residuals have been calculated and used in testing these data for the validity of the fundamental assumptions. Residuals are defined as plot yield plus twice the general mean minus the column mean minus the replication mean minus the variety mean.

We can roughly state the basic assumptions of the analysis-of-variance test in terms of residuals by saying that the residuals are normally distributed and that there is no pattern or system in the way in which they occur.

If we test by chi-square the hypothesis that the combined residuals of the 1939 and 1940 tests are normally distributed, then the resultant P = 0.06. Such a value of P indicates only a slight departure from normality. The distribution appears to be slightly peaked and positively skewed.

We should now examine the possibility that the residuals occur according to some plan or pattern.

If we consider that the residuals come at *random* from a normal population with a fixed standard deviation, then the six or seven residuals' occurring within a block (a block being the six or seven plots common to a given column and a given replication where the two cross at right angles) should be independent of the block-mean yield. These values are not independent for the 1939 and 1940 data. Thus when block-mean yields are correlated with block-mean residuals, the values of r for the Baart 38 strains in 1939 and 1940 are 0.23 and 0.24 respectively, and for the White Federation 38 strains for the same two years 0.23 and 0.54. Using the tables of David (6), one can calculate the probabilities of getting from a normal population for which the correlation is zero. correlation coefficients as high as the observed ones or higher. In samples of 25 these probabilities are 0.14, 0.13, 0.14, and 0.003 respectively. Let us compute, by the chi-square method of combining independent probabilities (6), the probability of the set of four observed values under the set of alternatives that the correlation coefficients of the sampled populations are unequal but all greater than zero. We find P = 0.0006. It is striking that the r values are the same for the three similar F values (see table 4), and that the much larger value of r occurs for the experiment with the exceedingly large F value.

Judging from the correlation between block-mean yields and block-mean residuals, too much has been subtracted from the poor plots and too little from the good plots in calculating the residuals. Hence, part of the variation in soil fertility has been assigned to the variation between varieties. That is, the design is too inflexible to take care of soil variation adequately from *one set of six or seven contiguous plots to another*. The result is a spreading effect on the strain means, and a tendency to indicate significance where none exists.

This correlation, furthermore, implies a parabolic relation between yield and the sum of squares of residuals.

Row totals (that is, a summation of yields of the plots occupying comparable

⁷ Six in the experiments with Baart 38 strains and seven for White Federation 38 strains.

positions across all replications) show pronounced coincident peaks of fertility culminating at the sixteenth row for both years for the Baart 38 strains, though the experiments occupied different areas in the two years. These coincident productivity peaks explain why, in the analysis of variance, the same Baart 38 strain appeared significantly higher in yield in both years. The row totals of White Federation 38 strains show a similar peak in 1939, but no very definite peak in 1940.

We may briefly summarize the evidence from these studies and from the previously mentioned work of Neyman, McCarthy, Salmon, and Baker relating directly to the assumptions on which the analysis of variance is based. According to both Neyman (10) and McCarthy (9) there are some cases of unequally correlated levels of fertility of plots assigned to different varieties, and such correlation causes overestimation of significance when the analysisof-variance technique is used. According to Baker (1), serious overestimation of significance may result from nonnormal distributions of yield estimates. According to Salmon (12) and others, the analysis of variance is invalid when the error variance differs from one part of the experiment to another. In the present work, correlation between fertility levels and residuals has been established. That is, the design has not prescribed a sufficiently flexible law of soilproductivity levels. This correlation means, not that residuals measured from their mean are more variable in one part of the experiment than in another, but that a bias in the residuals exists because soil productivity has been partially incorporated into strain differences. The bias to which we call attention is not attenuated as a cause of overestimation by the near identity of the strains tested. Certain conditions on soil-productivity levels, size of plot, and number of strains tested will lessen or make negligible the overestimation due to correlation between soil productivity and residuals.

APPLICATION TO OTHER YIELD-TESTING DESIGNS

So far as the authors are aware, no design now in use for testing large numbers of varieties is free from the danger of seriously overestimating significance when conditions are such that (1) fertility levels of plots assigned to different varieties are unequally correlated, (2) distribution of the estimates of a "true" plot yield is not normal, (3) the variances for different parts of the experiment are significantly different, and (4) an insufficiently flexible law is imposed on the productivity levels by an inadequate design. We believe, furthermore, that none of the conventional designs can adequately eliminate all these possible invalidating conditions under all conditions of testing. The lattice designs, now coming into vogue as the most efficient design for testing large numbers of varieties, have one admitted limitation : since varietal means are partially confounded with block effects, the use of these incomplete block designs may be cautioned against where large varietal differences are expected. In addition, the lattice designs seem to involve exactly the same difficulty experienced in these tests-namely, the danger of not subtracting the right amount from each plot yield or of not making the right "correction." They impose on the experiment a fixed formal framework, which may not be flexible enough to take care of spotted or abrupt changes in fertility levels.

Uniformity experiments are frequently recommended as preliminary steps

in determining the specific design best suited to test specific material in a given locality. Such tests, usually in operation for one year, or a very few years, cannot take into account the fluctuation of productivity levels from year to year as they are affected by changed environments, varying biological factors, tillage, and the like. They also provide no means of measuring the possible effect of complicating interactions when dissimilar rather than uniform material is under test.

We should not overlook the possibility that significance may be dangerously overestimated in any yield test of large numbers of varieties.

SUMMARY

An extensively used design, the modified Latin square, was used to test the comparative yielding ability of random selections from genetically similar component strains of Baart 38 and White Federation 38. Significant differences between strains were indicated in 1939. When the experiments were repeated in 1940, significant differences were again indicated, but with reversals from the previous year. When the results of the two years were combined, and the strains significantly greater and less than the general mean were again tested, these strains were found to be not significantly different. This last experiment was a small-scale Latin square experiment covering only a few strains. These facts prompted a critical examination of the data from the standpoint of the validity of the assumptions underlying the statistics used.

If the residuals in these experiments are examined, the block mean residual proves to be significantly and positively correlated with the block mean yield. Evidently, therefore, some of the difference in fertility levels of the plots has been assigned to strain differences. The result is a spreading effect on strain means, and a tendency to indicate significant differences in this large-scale experiment when, in fact, none exists.

Admittedly, small differences in yielding ability may actually exist between the strains tested in these experiments. Any such differences are masked, however, by the spotted variation in soil-fertility levels. Certainly the differences are not of the order of magnitude nor of the degree of significance indicated by the ordinary analysis of variance.

The demonstrated causes of overestimating significance in the analysis of variance are frequently assumed to be nonexistent. Indicated significant differences between "varieties" tested under conditions where any invalidating factors are operating should be viewed with skepticism. Feb., 1944]

LITERATURE CITED

1. BAKER, G. A.

- 1941. Fundamental distribution of errors for agricultural field trials. Natl. Math. Mag. 16:7-19.
- 2. BATEN, W. D., J. I. NORTHAM, and A. F. YEAGER.
 - 1941. Grouping of strains or varieties by use of a Latin square. Amer. Soc. Agron. Jour. 33:616-22.
- 3. BRIGGS, FRED N.
 - 1935. The backcross method in plant breeding. Amer. Soc. Agron. Jour. 27:971-73.

4. BRIGGS, FRED N.

- 1938. The use of the backcross in crop improvement. Amer. Nat. 72:285-92.
- 5. Cox, G. M., R. C. ECKHARDT, and W. G. COCHRAN.
 - 1939. The analysis of lattice and triple lattice experiments in corn varietal tests. Iowa Agr. Exp. Sta. Res. Bul. 281:1-66.
- 6. DAVID, F. N.
 - 1938. Tables of the correlation coefficient. Text 38 p. Tables 55 p. Biometrika Office, University College, London, Eng.
- 7. FISHER, R. A.
 - 1935. Discussion of: Yates, F. Complex experiments. Roy. Statis. Soc. Jour. Sup. 2:229-31.

8. JONES, DONALD FORSA.

1925. Genetics in plant and animal improvement. 568 p. John Wiley and Sons, Inc., New York, N. Y.

9. MCCARTHY, M. D.

1939. On the application of the Z-test to randomized blocks. Ann. Math. Statis. 10:337-59.

10. NEYMAN, J.

1935. Statistical problems in agricultural experimentation. Roy. Statis. Soc. Jour. Sup. 2:107-80.

11. Pope, O. A.

1936. Efficiency of single and double restrictions in randomized field trials with cotton when treated by the analysis of variance. Arkansas Agr. Exp. Sta. Bul. **326**:1-28.

12. SALMON, S. C.

1938. Generalized standard errors for evaluating bunt experiments with wheat. Amer. Soc. Agron. Jour. 30:647-63.

13. SNEDECOR, GEORGE W.

1934. Calculation and interpretation of variance and covariance. 96 p. Collegiate Press, Inc., Ames, Iowa.

14. SNEDECOR, GEORGE W.

1938. Statistical methods. rev. ed. 388 p. Collegiate Press, Inc., Ames, Iowa.

15. STRINGFIELD, G. H., R. D. LEWIS, and H. L. PFAFF.

1941. The 1940 Ohio coöperative corn performance tests. Ohio Agr. Exp. Sta. Spec. Cir. 61:1-30.

16. "Student."

1931. Yield trials. In: Hunter, H. Baillière's Encyclopedia of Scientific Agriculture. Vol. 2, p. 1342-61. Baillière, Tindall, and Cox, London, Eng.

- 17. SUNESON, C. A., O. C. RIDDLE, and F. N. BRIGGS.
 - 1941. Yields of varieties of wheat derived by backcrossing. Amer. Soc. Agron. Jour. 33:835-40.
- 18. WELCH, B. L.
 - 1937. On the Z-test in randomized blocks and Latin squares. Biometrika 29:21-52.

19. YATES, F.

1935. Complex experiments. Roy. Statis. Soc. Jour. Sup. 2:181-223, 243-47.

20. ZUBER, M. S.

1942. Relative efficiency of incomplete block designs using corn uniformity trial data. Amer. Soc. Agron. Jour. 34:30-47.

21. ZUBER, M. S., and J. L. ROBINSON.

,

1941. The 1940 Iowa corn yield test. Iowa Agr. Exp. Sta. Bul. n.s. P19:519-93.