

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

COMBINING EXPERIMENTS TO PREDICT FUTURE YIELD DATA¹

Foster B. Cady and David M. Allen²

ABSTRACT

Data from a series of fertility experiments including uncontrolled environmental variables are analyzed so that future yields may be predicted. A new criterion, the prediction sum of squares which is based on the performance of the estimated equation for predicting observations not included in the least squares estimation, is developed for selecting the best predictor variables. The procedure gives an equation with a minimal number of predictor variables and agronomically reasonable estimates of the regression coefficients.

Additional index words: data analysis, regression methodology, yield prediction, pooling of experiments.

Joint contribution of the Biometrics Unit, Plant Breeding and Biometry, Cornell University, Ithaca, New York, 14850, and the Department of Statistics, University of Kentucky. Paper No. BU-382-M of the Biometrics Unit.

² Professor of Biological Statistics, Cornell University and Assistant Professor of Statistics, University of Kentucky.

1 COMBINING EXPERIMENTS TO PREDICT FUTURE YIELD DATA¹2 Foster B. Cady and David M. Allen²

3
4 Inference to a population of yields, based on relatively few experi-
5 mental data points, is a key step in agronomic research. In methodology
6 it is stressed that the realized experiments should cover a range of
7 soil, management and climatic conditions, resulting in a series of experi-
8 ments over several locations and years. Voss et al, (5) emphasize the
9 importance of evaluating environmental factors in interpreting data from
10 a series of planned experiments. In the analysis of combined data, both
11 the observed yields from the experimentally controlled treatments and
12 measured data from the uncontrolled environmental factors are used to
13 estimate the relationship between the response variable and the various
14 hypothesized causal factors. This paper is chiefly concerned with the
15 development of a new procedure for the selection of variables in an
16 equation for predicting future yield data.

17
18 THE PREDICTION SUM OF SQUARES CRITERION

19 A major objective in the interpretation of data from a series of
20 experiments is the determination of a general yield equation. Suppose a

21 ¹ Joint contribution of the Biometrics Unit, Plant Breeding and Biometry,
22 Cornell University, Ithaca, New York, 14850, and the Department of
23 Statistics, University of Kentucky. Paper No. BU-382-M of the Biometrics
24 Unit.

25 ² Professor of Biological Statistics, Cornell University and Assistant
26 Professor of Statistics, University of Kentucky.

1 nitrogen fertility study with controlled levels of applied nitrogen is
2 conducted at each of several locations. A response model between yield
3 and applied nitrogen can be estimated at each site but an estimated model
4 over all sites is more important. Unfortunately, the statistical fitting
5 of all the data to one overall model is usually very poor, as exemplified
6 by Voss, et al (5). The same general function between, say, yield and
7 applied nitrogen, might exist for all sites but with different levels of
8 site variables, e.g., soil nitrogen, a different portion of the function
9 is observed at each site. Or a general yield function might be affected
10 through real interactions between the site variable and the experimentally
11 controlled variable. For either case, a number of uncontrollable factors
12 do exist and these cannot be held at constant or varying levels. However,
13 measurement can be made of these factors called site variables.

14 In a series of experiments designed for a combined analysis, a
15 relatively large number of site variables are measured, not only those
16 that the experimenter is highly certain are of value in prediction but
17 also questionable variables are initially included. The entire data set
18 is then used to indicate the variables to be included in the best esti-
19 mated yield equation. The addition of a variable to the prediction equa-
20 tion almost always increases (and never decreases) the variance of a pre-
21 dicted response, Walls and Weeks (6). However, failure to include an
22 important variable may result in a biased predicted observation. An
23 ideal procedure would select variables which are important in reducing
24 bias without selecting those which would unnecessarily add to the variance
25 of a predicted value.

26 The problem can now be stated as how one may select the predictor
27 variables in a general yield equation where some of the variables are

1 measured but not experimentally controlled, and are statistically cor-
2 related with each other. The commonly used selection procedures and
3 their drawbacks for use in combined analysis of fertility experiments
4 are discussed by Laird and Cady (4). These procedures, in general, use
5 the residual sum of squares as the criterion for variable selection; the
6 stepwise procedure will include a variable in the estimated model if a
7 statistical test for the reduction in the residual sum of squares is
8 significant at a predetermined type I error rate. The need for the
9 statistical test is to provide a stopping rule since the residual sum of
10 squares, in practice, will be lowered with each additional variable.

11 In the Laird and Cady (4) work, a criterion was needed to evaluate
12 the usefulness of three estimated reduced models. The commonly used
13 criteria of R^2 and the residual mean square were not sufficiently dis-
14 criminating. Since the estimated yield equations are used for predicting
15 future observations where sets of new values for the so-called independent
16 variables become available, it seemed reasonable to divide the data into
17 halves, estimate the parameters of the preselected alternative models with
18 one half, predict for the other half, and finally calculate the sum of
19 squares between the observed and predicted responses.

20 It then seemed reasonable to use a similar criterion in the develop-
21 ment of a new procedure for actually selecting the variables. This new
22 procedure would weigh the goodness of a particular potential variable to
23 predict observations not included in the estimation of the model para-
24 meters. The new criterion, the prediction sum of squares (PRESS), is
25 based on the difference between Y_i , the observed response and $\hat{Y}_{(i)}$, the
26 response predicted from an estimated equation excluding the i^{th} observa-
27 tion, i.e., the predicted response is based on $n-1$ observations. The

1 deviations are then squared and added over the n observations, i.e.

$$2 \text{ PRESS} = \sum_{i=1}^n (y_i - \hat{y}_{(i)})^2 .$$

3
4 The sequential PRESS algorithm (SPA), utilizing the new criterion,
5 is a procedure for selecting the predictor variables. SPA is analagous
6 to the stepwise procedure using the residual sum of squares as the
7 criterion. The one variable having the smallest value of PRESS is
8 entered first. At each subsequent stage, the variable having the small-
9 est value of PRESS in conjunction with the previously entered variable(s)
10 is entered. For example, at the second stage, that variable among the
11 remaining potential variables giving the largest reduction in the
12 prediction sum of squares is selected. Allen (1) gives an efficient
13 algorithm for this sequential procedure.

14 When selecting variables by SPA, PRESS typically decreases rapidly
15 with the first few variables selected. Each of the next few variables
16 decreases PRESS by small amounts until a minimum is obtained, after which
17 additional variables increase PRESS. The net effect, when viewing a
18 plot of PRESS with respect to the number of variables in the equation
19 is a backward J shaped curve with a nearly flat bottom. Continuing to
20 the minimum appears to reduce the biases squared relative to the variances,
21 but the variances are known to be increasing at each stage and the biases
22 have to be estimated, suggesting that the cut-off point could be made
23 subjectively before the minimum is observed.

24 RESULTS

25
26 The procedure SPA was compared with stepwise regression using data
27 from a series of 76 experiments in the Bajio area of Mexico that were

1 carried out over a four year period and described by Laird and Cady (4).
2 Four experiments were deleted in order to have a total number divisible
3 into groups, a facility not used in the results reported here. The full
4 model included the 33 independent variables listed in Table 1 with the
5 estimated partial regression coefficients. Stepwise regression using a
6 type I error rate of 5% as a stopping rule was used and the resulting
7 estimated partial regression coefficients of the 17 entering variables are
8 also given in Table 1. With SPA the variables entered in the order as
9 shown in Figure 1. The prediction sum of squares drops rapidly with the
10 first few variables to enter the equation. The figure does not show the
11 minimum value of the prediction sum of squares of 116.668 nor the follow-
12 ing increase to 122.315 with all variables entered. The type of function
13 would indicate that a cut-off for stopping variables entering the equa-
14 tion could be made before the minimum is reached. In this example the
15 cut-off was made after the soil slope variable entered and the estimated
16 partial regression coefficients are given in Table 1. The constant or
17 mean was the single variable that most reduced the prediction sum of
18 squares at the first stage but theoretically it could have been another
19 variable, as exemplified later.

20 In practice SPA is used to select predictor variables. In this situa-
21 tion the experimenter would like to determine those predictor variables
22 that will predict well future responses when new sets of independent vari-
23 able values become available. The Bajic data from the first three years
24 (228 observations) were used for selecting predictor variables by both
25 stepwise and SPA. Using a 5% significance level for a stopping rule, 15
26 of the potential 33 predictor values were selected by stepwise, including
27 K, N, N², CA, CB, D, DN, DB, F, G, G², B², BA, H, J, and LA. The

1 Table 1. The independent variables (the ranges are given in parentheses)
 2 and the estimated partial regression coefficients for the full and reduced
 3 models.

4	Independent Variable	Symbol	Estimated Regression Coefficients		
			Full	Stepwise	SPA
5	Constant	K	-0.3170	-0.5446	1.5780
6	Applied nitrogen (linear) (0-3)	N	1.8410	1.8050	1.4540
7	Applied nitrogen (quadratic)	N ²	-0.1552	-0.1547	-0.1528
8	Total soil nitrogen (linear) (46-137)	A	-0.0290		0.0098
9	Total soil nitrogen (quadratic)	A ²	0.0150	0.0032	
10	A x N	AN	-0.0396	-0.0406	
11	Previous crop (linear) (10-25)	B	0.2220		
12	Previous crop (quadratic)	B ²	-0.0813	-0.0176	
13	B x N	BN	-0.0014		
14	B x A	BA	0.0771	0.0711	
15	Excess moisture (0-7)	C	0.1066	-0.2656	-0.2436
16	C x N	CN	-0.0374		
17	C x A	CA	-0.0217		
18	C x B	CB	-0.0794		
19	Drought (0-112)	D	0.0309		
20	D x N	DN	-0.0096	-0.0091	-0.0091
21	D x A	DA	-0.0023		
22	D x B	DB	-0.0259		
23	Depth of rooting zone (25-99)	E	-0.0054		
24	Soil slope (1-80)	F	-0.0124	-0.0111	-0.0086
25	Soil texture (linear)(2-5)	G	1.2800	1.2740	
26	Soil texture (quadratic)	G ²	-0.1591	-0.1630	
27	Hail (0-5)	H	0.5556	0.2651	-0.2737
28	H x N	HN	-0.0003		
29	H x A	HA	-0.0802	-0.0694	
30	H x B	HB	-0.0159		
31	Blight (H. Turcicum)(0-8)	J	-1.089	-0.2733	-0.2677
32	J x N	JN	0.0183		
33	J x A	JA	0.0611		
34	J x B	JB	0.0139		
35	Weeds (0-4)	L	-1.7750	-0.9231	
36	L x N	LN	-0.0004		
37	L x A	LA	0.1111	0.0757	
38	L x B	LB	0.0458	0.0183	

1 prediction sum of squares for SPA is shown in Figure 2. Again, after the
2 first few variables, the curve is nearly flat as the minimum of 88.099
3 (not shown in the figure) is approached. With all variables included the
4 prediction sum of squares is 93.819. As before, the cut-off was made
5 after the soil slope variable, the linear effect of applied nitrogen
6 giving the smallest prediction sum of squares at the first stage.

7 Rather surprising results from using these estimated models, includ-
8 ing the full model, are shown in Table 2. Though the residual mean square
9 using SPA is 10% higher than the stepwise, the future prediction sum of
10 squares (as defined in the table) with SPA selected predictor variables
11 is approximately 30% lower than the stepwise. As expected, the full
12 model fits the data from the first three years better than either reduced
13 estimated model as evidenced by the smallest residual mean square. How-
14 ever, the predictive ability of the full model is poor, more than doubling
15 in the future prediction sum of squares when compared with the SPA esti-
16 mated model.

17 DISCUSSION

18 The agronomist prefers a general yield equation with regression co-
19 efficients that are meaningful in both a qualitative and quantitative sense.
20 The problems in interpreting estimated coefficients with correlated vari-
21 ables have been discussed by Voss et al, (5). As shown in Table 1, and
22 conjectured in general, SPA selected variables yield estimates of the
23 coefficients with proper signs and reasonable magnitudes. For example,
24 SPA selected the linear effect of soil nitrogen; the estimated coefficient
25 is positive, giving nearly a ton increase in yield for a field reasonably
26 high in soil nitrogen. However, the stepwise procedure selected the quad-
27 ratic effect of soil nitrogen and the interaction with applied nitrogen
with estimates that cannot be interpreted quantitatively without

1
2 Table 2. Results from using the first three years of data (228 observa-
3 tions) for predictor variable selection and subsequently predicting for
4 the fourth year (60 observations).

	Estimated Model		
	Full	Stepwise	SPA
7 Number of predictor variables	33	15	9
8 Residual sum of squares 9 (228 observations)	67.521	80.134	91.806
10 Residual mean square	0.348	0.383	0.421
11 Prediction sum of squares (228 observations)	93.819	92.868	100.794
12 Future prediction sum of squares†	67.661	42.395	30.846

13
14
15
16
17 † Sixty squared deviations between the predicted responses, using the
18 first three years of data for the predictor variable selection, and the
19 observed fourth year responses summed over the fourth year.
20
21
22
23
24
25
26
27

1 considering the constant term and other variables. Using the estimated
2 model selected by SPA, yield may be plotted as a quadratic polynomial
3 function of applied nitrogen. Since the drought by applied nitrogen
4 interaction is important, the plot would have a different estimated
5 function for each of several drought levels. These curves would then be
6 influenced by the main effects of soil nitrogen, excess moisture, slope,
7 hail, and blight. Such a plot is Figure 3 using average values of soil
8 nitrogen and slope, low values for hail and blight and zero for excess
9 moisture.

10 As demonstrated, the SPA procedure can be used for analyzing and
11 interpreting an existing set of data. However, a more important role of
12 SPA is the selection of predictor variables in situations where the
13 estimated model will be used for predicting yields at a later date with
14 new sets of independent variable values. It is assumed that the new
15 values are in the same range as those used in the predictor variable
16 selection stage. For example, with soil test recommendations, a series
17 of experiments can be conducted in a given area over a period of time.
18 Based on these experiments, a prediction equation is desired which can be
19 used in future years with soil test values, knowledge of probabilities
20 involving climatic variables and farmer information on management vari-
21 ables. Now the emphasis in the variable selection procedure primarily is
22 the selection of those variables which will do well in predicting future
23 observations; the secondary emphasis is the selection of those variables
24 which will do the best in reducing the residual sum of squares for that
25 particular set of data, which is important if the interpretation of the
26 existing set of data is the sole objective.

27 Most of the variables selected by SPA were also selected by the

1 stepwise procedure. In fact, the ordering of the first few variables is
2 the same for the two procedures and curves for the residual sum of squares
3 analogous to those in Figures 1 and 2 are similar for the first few
4 variables. However, the residual sum of squares will in practice continue
5 to decrease as more variables enter the equation. In addition the dif-
6 ferences between the two procedures in the variables entering during the
7 earlier stages can be important. The linear effect of soil nitrogen is
8 an example in the data set presented here. The two procedures were also
9 compared with data from Voss et al, (5) with 29 potential predictor
10 variables. SPA selected the linear effect of plant population, a vari-
11 able not selected by the stepwise procedure which selected several inter-
12 actions in the prediction model. Of the ten variables selected by PRESS,
13 nine were linear effects of site variables, the other being the inter-
14 action between drought and applied nitrogen.

15 A recommendation could be made for using the stepwise procedure but
16 with an earlier cut-off point, e.g., using smaller significance levels
17 than .05 or perhaps not using any hypothesis testing. It is believed
18 that the stepwise procedure as described by referenced textbooks, e.g.
19 Draper and Smith (3), historically is too deeply associated with hypo-
20 thesis testing to make any real break with standard use. Also, it is
21 believed that the prevailing practice is to use the larger significance
22 levels in variable selection so that a so-called important variable
23 would not be eliminated from the equation. This practice for including
24 a large number of variables in an equation is based, in part, on the
25 feeling that extra variables can't do any harm. With different approaches,
26 it has been shown in both the Laird and Cady (4) and the present study
27 that extra variables can indeed be detrimental for prediction purposes.

1 The SPA procedure, based on the prediction sum of squares criterion, is
2 not associated with any hypothesis testing. The cut-off point is deter-
3 mined by the minimum value of the criterion. With small data sets, the
4 minimum is rapidly reached as variables are entered and the increase in
5 the criterion immediately begins after the minimum is reached. An
6 example is given by Allen (1). With large data sets, such as the one
7 used in this study, a practical recommendation is to use a cut-off before
8 the absolute minimum is reached. This is arbitrary, of course, but the
9 decision is not difficult with the rapid decrease in the prediction sum
10 of squares before reaching the flat portion of the curve.

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

REFERENCES

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27

1. Allen, David M. 1971. The prediction sum of squares as a criterion for selecting predictor variables. Submitted to Technometrics.
2. Colwell, J. D. and K. M. Stackhouse. 1970. Some problems in the estimation of simultaneous fertilizer requirements of crops from response surfaces. Australian J. of Exp. Agr. and Anim. Husb. 10: 183-195.
3. Draper, N. R. and H. Smith. 1966. Applied regression analysis. John Wiley and Sons, Inc., New York.
4. Laird, R. J. and F. B. Cady. 1969. Combined analysis of yield data from fertilizer experiments. Agron. J. 61:829-834.
5. Voss, R. E., J. J. Hanway and W. A. Fuller. 1970. Influence of soil, management, and climatic factors on the yield response by corn (*Zea Mays* L.) to N, P and K fertilizer. Agron. J. 62:736-740.
6. Walls, R. C. and D. L. Weeks. 1969. A note on the variance of a predicted response in regression. Amer. Statist. 23:24-26.

