DOCUMENT RESUME

ED 121 603

AUTHOR Levin, Joel R.; Subkoviak, Michael J.

TITLE Additional Considerations in Determining Sample

Size.

INSTITUTION Wisconsin Univ., Madison. Dept. of Educational

Psychology.

REPORT NO Occas-Pap-10

PUB DATE 75

NOTE 17p.; Paper presented at the Annual Meeting of the

American Educational Research Association (San

SE 020 659

Francisco, California, April 19-23, 1976)

EDRS PRICE MF-\$0.83 HC-\$1.67 Plus Postage

DESCRIPTORS Analysis of Variance; *Data Analysis; *Data

Collection: *Educational Research: *Research;

*Sampling; Statistical Analysis

ABSTRACT

Levin's (1975) sample-size determination procedure for completely randomized analysis of variance designs is extended to designs in which antecedent or blocking variables information is considered. In particular, a researcher's choice of designs is framed in terms of determining the respective sample sizes necessary to detect specified contrasts of a given magnitude with given Type I and Type II errors. A solution is provided for dealing with real-world considerations in which errors of measurement cannot be neglected. A worked example presents an instance wherein a blocking strategy is clearly advantageous assuming infallible measuring instruments, but not when the same instruments are granted fallibility. (Author/SD)



Session No. 13.15 (under a different title)

U S DEPARTMENT OF HEALTH, EDUCATION & WELFARE NATIONAL INSTITUTE OF EDUCATION

THIS DOCUMENT HAS BEEN REPRO-DUCED EXACTLY AS RECEIVED FROM THE PERSON OR ORGANIZATION ORIGIN-ATING IT POINTS OF VIEW OR OPINIONS STATED DO NOT NECESSARILY REPRE-SENT OFFICIAL NATIONAL INSTITUTE OF EDUCATION POSITION OR POLICY

5 B

ADDITIONAL CONSIDERATIONS IN DETERMINING

SAMPLE SIZE

Joel R. Levin

and

Michael J. Subkoviak

Occasional Paper No. 10
Laboratory of Experimental Design
Department of Educational Psychology
University of Wisconsin

1975

Additional Considerations in Determining
Sample Size

Joel R. Levin and Michael J. Subkoviak
University of Wisconsin

ABSTRACT

Levin's (1975) sample-size determination procedure for completely randomized analysis of variance designs is extended to designs in which antecedent or blocking variable information is considered. In particular, a researcher's choice of designs is framed in terms of determining the respective sample sizes necessary to detect specified contrasts of a given magnitude with given Type I and Type II errors. A solution is provided for dealing with real-world considerations in which errors of measurement cannot be neglected. A worked example presents an instance wherein a blocking strategy is clearly advantageous assuming infallible measuring instruments, but not when the same instruments are granted fallibility.



Additional Considerations in Determining Sample Size

Joel R. Levin and Michael J. Subkoviak University of Wisconsin

INTRODUCTION

When it comes to designing an experiment, an educational researcher can draw from a variety of sources--some in the form of old wives' tales, and some in the form of theoretically sound recommendations (e.g., Feldt, 1958)--to determine whether it is preferable to assign subjects randomly to K experimental conditions and subsequently to perform an analysis of variance on the dependent variable Y (hereafter referred to as a completely randomized design); or rather to include in the analysis antecedent information based on variable X (known or assumed to be related to Y). The antecedent information included can be operationally dealt with in various ways: chiefly, in terms of randomized blocks analysis, analysis of covariance, or analysis of an index of response (such as change scores)--cf. Porter & Chibucos (1974).

The major advantage of these procedures, relative to the completely randomized design, is one of reducing the within-treatment variability by removing the variation in Y that is due to the relationship between X and Y. The present paper focuses on one of these procedures, namely the randomized block design, as a competitor to the completely randomized design; and, in particular, it considers an alternative to the traditional way of deciding whether to block or not to block that includes real-world situations in which errors of measure-

ment associated with X, Y, or both are likely to be present. Moreover, since the discussion by Porter and Chibucos (1974) suggests that in "true" (Campbell & Stanley, 1966) experiments of moderate sample size, analysis of covariance and analysis of an index of response may be regarded as essentially equivalent procedures to blocking--within degrees-of-freedom differences and slight differences in their error expected mean squares--the material presented here has implications for the other two procedures as well.

Reliability and Sample Size

Statistics texts typically acknowledge four ingredients of hypothesis testing: (a) Type I error probability (α); (b) Type II error probability (β) or its complement, power (1 - β); (c) sample size, and (d) the magnitude of the experimental effect of interest. In planning an experiment, a researcher can specify α and the power desired to detect an effect of specified magnitude, and subsequently calculate the required sample size; or, in evaluating a completed experiment, the predetermined α level and sample size can be used to compute the power that was available to detect an effect of given magnitude.

Such calculations tacitly assume that dependent variables and/or antecedent variables are measured without error, i.e., they are perfectly reliable (true scores). In actual practice, however, both antecedent and dependent variables are measured with error, i.e., they are fallible (observed scores), with the result that "textbook" power/sample size calculations (which do not take the unreliability of the observed data into account) produce inaccurate estimates. In particular, they produce underestimates of required sample sizes in the planning stage and overestimates of available power in the post hoc evaluation condition. The present paper provides formulas for the computation of power and sample size that include the reliability coefficient of observed scores, thereby augmenting the list of hypothesis-testing ingredients mentioned above.



Several authors have considered the effect of unreliability on statistical tests (e.g., Cleary & Linn, 1969; Cleary, Linn, & Walster, 1970; Overall & Dalal, 1965; Sutcliffe, 1958; Porter, Note 1). Cleary et al. (1970), for example, have demonstrated that the power of the F-test in a one-way, fixedeffects analysis of variance (ANOVA) decreases as the reliability--and also as the validity--of the dependent variable decreases. The purpose of the present paper is to extend some of the Cleary et al. notions to designs in which antecedent information is considered; in particular, to the randomized block design. Moreover, in contrast to the commonly recommended strategy for deciding whether or not it would be advantageous to block (i.e., by determining the relative efficiency of a randomized block design to a completely randomized design for a fixed number of subjects--cf. Kirk, 1968, pp. 147-149), the strategy adopted here consists of framing the decision in terms of the respective sample sizes associated with the two designs that are required to yield equivalent power for detecting specified effects of interest (see, for example, Cohen, 1969, pp. 46-50).

CASE 1: LATENT TRUE VARIABLES

Sample Size Determination for the Completely Randomized Design

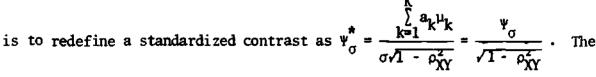
The reader is referred to Levin (1975) for a discussion of sample size determination based on a researcher's <u>a priori</u> specification of the minimum value of any given linear <u>contrast</u> of interest (which has been called Ψ_{σ}) in accordance with desired α and 1- β . The resulting number of subjects required per experimental condition (n) guarantees the researcher the desired power associated with detecting the contrast of interest, should it be of the specified magnitude. In the case of a planned-comparison approach to hypothesis testing, the F-test is performed with 1 and K(n - 1) degrees of freedom (these



referring to the degrees of freedom associated with the contrast and the mean square within respectively); and in this situation the probability of detecting a contrast of the magnitude specified is alternatively the probability of obtaining a significant \underline{F} -ratio (both 1 - β). In the case of a post hoc approach to hypothesis testing, the \underline{F} -test is performed with K - 1 and K(n - 1) degrees of freedom (where K - 1 represents the degrees of freedom associated with the mean square between); and in this situation the probability of detecting a contrast of the magnitude specified is alternatively the probability of obtaining a significant \underline{F} -ratio and then identifying that contrast as statistically significant according to Scheffé's (1953) multiple comparison procedure (see Levin, 1975). According to this formulation, Ψ_{σ} represents the magnitude of the contrast in means considered to be of interest to the researcher, and which is expressed in within-treatment standard deviation units (σ). Thus, if $\Psi = \sum_{k=1}^K a_k \mu_k$ (where the a_k represent contrast coefficients chosen such that $\sum_{k=1}^K a_k \mu_k = 0$), then $\Psi_{\sigma} = \frac{\sum_{k=1}^K a_k \mu_k}{\sigma}$.

Sample Size Determination for the Randomized Block Design

Rather than adopting the completely randomized design, a researcher may choose to form n blocks of K subjects (on the basis of some relevant antecedent information), and then randomly assign subjects within blocks to the K treatment conditions. It is well known that the effect of introducing a blocking variable into the design is to reduce σ by a factor of $\sqrt{1-\rho_{XY}^2}$, where ρ_{XY} represents the correlation between the antecedent variable and the dependent variable. Thus, in terms of the present approach, all that needs to be done





effect of blocking, then is to increase the value of Ψ_{σ} of the completely randomized design which, if it overcompensates for the corresponding loss in error degrees of freedom, i.e., from K(n-1) to (K-1)(n-1), results in a decrease in the number of subjects required in order to maintain equivalent power to that in the completely randomized case.

CASE 2: FALLIBLE VARIABLES

The above discussion has proceeded under the assumption that the only "error" in the ANOVA model consists of <u>subject</u> error. If there is <u>measurement</u> error as well, one's <u>effective</u> power will not be as great as one's <u>nominal</u> power; or, stated differently, a researcher will require more subjects than the "textbook" sample size determination indicates are needed in order to have the desired power (see, for example, Cleary, et al., 1970). Classical test theory (Lord & Novick, 1968) assumes that the observed score Y_i for person i is equal to his or her true score T_i plus measurement error E_i , such that $Y_i = T_i + E_i$. Since T_i and E_i are independently distributed with respective expected values of μ_T and 0 and respective variances of σ_T^2 and σ_E^2 , it follows that:

$$\mu_{\Upsilon} = \mu_{\tilde{T}} + 0$$

$$= \mu_{\tilde{T}}$$
(1)

and

$$\sigma_{\rm Y}^2 = \sigma_{\rm T}^2 + \sigma_{\rm E}^2 \tag{2}$$

The reliability of observed scores Y_i is the ratio of true score variance to observed score variance:

$$\rho_{YY}^{1} = \frac{\sigma_{T}^{2}}{\sigma_{T}^{2} + \sigma_{F}^{2}} = \frac{\sigma_{T}^{2}}{\sigma_{Y}^{2}}$$
(3)



Sample Size Determination for the Completely Randomized Design

randomized design? As was noted previously, Ψ_{σ} is simply a contrast involving the treatment means which is expressed in within-treatment standard deviation $\sum_{k=1}^{K} a_k u_k$ units, or $\Psi_{\sigma} = \frac{k=1}{\sigma}$. Because of the relationship in (1), the numerator of Ψ_{σ} is unaffected by measurement errors. What is affected is the denominator. Thus, σ in Ψ_{σ} reflects the within-treatment standard deviation of true scores, or σ_T . Following Cleary et al. (1970) and employing (3), we note that in terms of observed scores, $\sigma_Y = \frac{\sigma_T}{\sqrt{\rho_{YY}}}$. Thus, for the usual case where measurement errors associated with the dependent variable are expected, we simply redefine

How do these properties affect sample size determination in the completely

$$\underline{\Psi}_{\sigma} = \frac{\sum_{k=1}^{K} a_k \mu_k}{\sigma_{T/\sqrt{\rho_{YY}}}} = \sqrt{\rho_{YY}}, \ \Psi_{\sigma}$$

where it may be easily shown (though it will not be here) that ρ_{YY} , represents the (assumed common) within-treatment reliability of the dependent variable.

Sample Size Determination for the Randomized Block Design

In the case of the randomized block design, the situation becomes complicated due to potential errors of measurement associated with X in addition to those associated with Y. Employing correction-for-attenuation formulas, one can obtain the following general expression:

$$\underline{\Psi}_{\sigma}^{*} = \sqrt{\frac{\rho_{XX}, \rho_{YY}, -\rho_{XY}^{2}}{\rho_{XX}, (1-\rho_{XY}^{2})}} \Psi_{\sigma}^{*}$$

(where $\rho_{\chi\chi'}$ represents the reliability of the antecedent variable).



Ψ_σ as:

It should be noted that this expression can be easily adapted to fit various special cases. In particular, if only X is assumed to be fallible, it may be seen that:

$$\underline{\underline{\Psi}}_{\sigma}^{*} = \sqrt{\frac{\rho_{XX'} - \rho_{XY}^{2}}{\rho_{XX'}(1 - \rho_{XY}^{2})}} \, \underline{\Psi}_{\sigma}^{*}.$$

On the other hand, if only Y is fallible:

$$\underline{\underline{\Psi}}_{\sigma}^{*} = \sqrt{\frac{\rho_{YY}^{*} - \rho_{XY}^{2}}{1 - \rho_{XY}^{2}}} \, \underline{\Psi}_{\sigma}^{*}.$$

Finally, if neither X nor Y is fallible:

$$\underline{\underline{\Psi}}_{\sigma}^{*} = \sqrt{\frac{1 - \rho_{XY}^{2}}{1 - \rho_{XY}^{2}}} \, \underline{\Psi}_{\sigma}^{*} = \underline{\Psi}_{\sigma}^{*}$$

which is as it should be.

AN EXAMPLE

Levin's (1975) sample size determination formula is given by:

$$\phi = \sqrt{\frac{n\Psi_{\sigma}^{2}}{(\nu_{1}+1)\sum_{k=1}^{K} a_{k}^{2}}}$$
 (4)

where: ϕ = a parameter in the Pearson and Hartley (1951) power charts, available in most experimental design textbooks; more complete tables displaying ϕ are also available (e.g., Tiku, 1967, 1972).

Let us apply (4) to the simplest ANOVA situation, namely for K = 2 which is equivalent to the independent two-sample (nondirectional) t-test situation.



Assume that a researcher wishes to have an 80 percent chance of detecting a difference in K = 2 means of at least 1 standard deviation unit, based on a Type I error probability of .05. How many subjects per treatment group should he/she include? [With reference to Formula (4), it should be noted that for all cases to be considered, $v_1 + 1 = 2$, which will always equal K in the one-way layout; and $\sum_{k=1}^{K} a_k^2 = 2$, which will always be true when only pairwise differences in means are of interest, even for K > 2. (However, in some situations complex comparisons may interest the researcher, in which case the value of $\sum_{k=1}^{K} a_k^2$ will change--see Levin, 1975.)]

The information contained in the preceding paragraph may be translated as follows: α = .05, 1 - β = .80, Ψ_{σ} = 1.00. Incorporating this into (4) and the appropriate power charts, and proceeding in the manner described by Levin, we find that in the <u>completely randomized</u> situation (assuming a <u>perfectly reliable dependent variable</u>), a total of 17 subjects per treatment group is required to yield the desired power.

If we further assume that an antecedent variable is selected that correlates .50 with performance on the dependent measure (i.e., ρ_{XY} = .50), then it can be seen that $\Psi_{\sigma}^{*} = \frac{1}{\sqrt{1-(.50)^{2}}}$ = 1.155. Substituting this into (4) and checking with the appropriate ν_{2} , we find that if a randomized block design (assuming perfectly reliable antecedent and dependent variables) were employed, a total of 14 subjects per treatment group would be required to yield equivalent power to that in the completely randomized design above.

Now let us suppose that either or both of the two variables involved (antecedent and dependent) are fallible. Given separate (and equal) reliabilities of $\rho_{\chi\chi^1} = \rho_{\gamma\gamma^1} = .80$, for example, we are able to retrace the steps associated with (4), incorporating $\underline{\Psi}_{\sigma}$ and $\underline{\underline{\Psi}}_{\sigma}^*$ as previously defined. Table 1 summarizes

Insert Table 1 about here

the results of this endeavor.

What is especially interesting about this particular example is that even though we start out with a situation in which it is clearly preferable to block (as reflected by a total savings of six subjects for Situation 1 of Table 1), by the time the antecedent and dependent variables are both granted fallibility on the order of $\rho_{\chi\chi'} = \rho_{\gamma\gamma'} = .80$, the randomized block advantage disappears (as reflected by the 0 total subject savings difference in Situation 4 of Table 1).

To make this lesson somewhat more concrete, assume that a researcher is interested in comparing the efficacy of two instructional variations designed to teach eighth grade mathematics. Both variations are to be incorporated into programed instruction booklets and randomly assigned to students within classrooms or schools), and end-of-year performance will be assessed via a standardized mathematics achievement test. Suppose further in this hypothetical situation that the production cost of the booklets is somewhat of a factor, so that an experimental design that will yield the desired power with the fewest students is the one to be selected. Given this information, should the researcher randomly assign students to the two treatment conditions or block on seventh grade standardized mathematics achievement scores, known to correlated .50 with eighth grade scores? Ignoring the unreliability associated with two achievement tests (as in the "textbook' case), the researcher would clearly do well to block; he would require six fewer students with a randomized block design than with a completely randomized design. However, considering the published reliabilities of the two tests of .80, the researcher would discover that it makes little difference which of the two experimental designs he selects, since there is a 0 subject savings. In fact, if it would require some additional effort to obtain and/or record the seventh grade achievement data the researcher may well



opt for the seemingly less efficient (though not so in this case) completely randomized design.

CONCLUSION

This particular example is but one of several that could have been contrived. What should be clear to the reader, based on this example and the larger message of this paper, is as follows: First, each potential experiment should be examined on an <u>a priori</u> basis to determine whether or not it is advantageous to block. This decision cannot be made without considering the number of treatment conditions included, the magnitude of the relationship between the antecedent and blocking variables (ρ_{XY}) , as well as the various hypothesis-testing ingredients described at the outset of the paper. Second, to follow these procedures without simultaneously considering errors of measurement is to live in a "fool's paradise," for these too will affect block-no block decisions. In cases where <u>a priori</u> reliability information is lacking, pilot research or sagacious judgments (to obtain approximate and conservative estimates, respectively) will surely do better than nothing.



Reference Note

Porter, A. C. <u>The effects of using fallible variables in the analysis of covariance</u>. Unpublished doctoral dissertation, University of Wisconsin, Madison, 1967.



References

- Campbell, D. T., & Stanley, J. C. Experimental and quasi-experimental designs for research. Chicago: Rand McNally, 1966.
- Cleary, T. A., & Linn, R. L. Error of measurement and the power of a statistical test. British Journal of Mathematical and Statistical Psychology, 1969, 22, 49-55.
- Cleary, T. A., Linn, R. L., & Walster, G. W. Effect of reliability and validity on power of statistical tests. In E. F. Borgatta and G. W. Bohrnstedt (Eds.), Sociological methodology, San Francisco: Jossey-Bass, 1970.
- Cohen, J. Statistical power analysis for the behavioral sciences. New York:

 Academic Press, 1969.
- Feldt, L. S. A comparison of the precision of three experimental designs employing a concomitant variable. <u>Psychometrika</u>, 1958, <u>23</u>, 335-353.
- Kirk, R. E. <u>Experimental design: Procedures for the behavioral sciences</u>. Bel-mont, California: Brooks/Cole, 1968.
- Levin, J. R. Determining sample size for planned and post hoc analysis of variance comparisons. <u>Journal of Educational Measurement</u>, 1975, <u>12</u>, 99-108.
- Lord, F. M., & Novick, M. R. <u>Statistical theories of mental test scores</u>. Reading, Mass.: Addison-Wesley, 1968.
- Overall, J. E., & Dalal, S. N. Design of experiments to maximize power relative to cost. Psychological Bulletin, 1965, 64, 339-350.
- Pearson, E. S., & Hartley, H. O. Charts of the power function for analysis of variance tests, derived from the non-central F-distribution. <u>Biometrika</u>, 1951, 38, 112-130.
- Porter, A. C., & Chibucos, T. R. Analysis issues in summative evaluation. In G. Borich (Ed.), Evaluating educational programs and products. Englewood Cliffs, N. J.: Educational Technology Press, 1974.



- Scheffé, H. A method for judging all contrasts in the analysis of variance.

 Biometrika, 1953, 40, 87-104.
- Sutcliffe, J. P. Error of measurement and the sensitivity of a test of significance. <u>Psychometrika</u>, 1958, <u>23</u>, 9-17.
- Tiku, M. L. Tables of the power of the F-test. <u>Journal of the American</u>
 <u>Statistical Association</u>, 1967, 62, 525-539.
- Tiku, M. L. More tables of the power of the <u>F</u>-test. <u>Journal of the American</u>
 <u>Statistical Association</u>, 1972, <u>67</u>, 709-710.

Table 1. Comparison of Completely Randomized (CR) and Randomized Block (RB) Design Sample Sizes for the Present Example $(K=2,\;\alpha=.05,\;1-\beta=.80,\;\Psi_{_{\hbox{\scriptsize C}}}=1.00,\;\rho_{\hbox{\scriptsize XY}}=.05)$

	Situation	Ψ _σ or Equivalent	Number of Subjects Per Group	Total Subject Savings (RB - CR)
1.	X is Infallible,			
	Y is Infallible			
	CR	1.000	17	6
r	RB	1.155	14	
2.	X is Infallible,			
	Y is Fallible (ρ_{YY} , = .80)			
17	CR	.894	21	6
•	RB	.989	18	
3.	X is Fallible (ρ_{XX} , = .80),			
	Y is Infallible			
	CR	1.000	17	4
	RB	1.106	15	
4.	X is Fallible ($\rho_{\chi\chi}$, = .80),			
	Y is Fallible (ρ_{YY} , = .80)			
	CR	.894	21	n
C [*]	RB	.931	21	