

DOCUMENT RESUME

ED 476 918

TM 034 969

AUTHOR Newman, Isadore; McNeil, Keith; Fraas, John W.
TITLE Deja Vu: Another Call for Replications of Research, Again.
PUB DATE 2003-04-00
NOTE 16p.; Paper presented at the Annual Meeting of the American Educational Research Association (Chicago, IL, April 21-25, 2003).
PUB TYPE Reports - Descriptive (141) -- Speeches/Meeting Papers (150)
EDRS PRICE EDRS Price MF01/PC01 Plus Postage.
DESCRIPTORS *Effect Size; *Research Methodology; *Statistical Significance
IDENTIFIERS *Research Replication

ABSTRACT

Over the last few years, there has been evolution, although not a linear one, that has progressed from an emphasis on statistical significance to an emphasis on effect size to an emphasis on both of these concepts to what is believed to be a pragmatic emphasis on replicability. This paper presented two methods of estimating a study's replicability that researchers should consider reporting along with their statistically significant and effect size findings. One method of estimating the replicability of the findings deals with replication in the exact same system. The second method, which may contain subjective probability values, is used to estimate the replicability of a study's findings in a system that may differ from the initial system with respect to salient variables. The incorporation of the replicability estimates delineated in this paper would provide critical information to decision makers about the likelihood that the implementation of a particular method or treatment would produce similar results in their systems. (Contains 2 tables and 12 references.) (Author/SLD)

ED 476 918

Déjà vu: Another Call for Replications of Research, Again

Isadore Newman
University of Akron

Keith McNeil
New Mexico State University

John W. Fraas
Ashland University

Paper presented at the annual meeting of the American Educational Research Association,
Chicago, IL, April 21-25, 2003

TM034969

U.S. DEPARTMENT OF EDUCATION
Office of Educational Research and Improvement
EDUCATIONAL RESOURCES INFORMATION
CENTER (ERIC)

This document has been reproduced as
received from the person or organization
originating it.

Minor changes have been made to
improve reproduction quality.

• Points of view or opinions stated in this
document do not necessarily represent
official OERI position or policy.

PERMISSION TO REPRODUCE AND
DISSEMINATE THIS MATERIAL HAS
BEEN GRANTED BY

J. W. Fraas

TO THE EDUCATIONAL RESOURCES
INFORMATION CENTER (ERIC)

1

Abstract

Over the last few years there has been an evolution, not a linear one however, that has progressed from an emphasis of statistical significance to an emphasis on effect size to an emphasis on both of these concepts to what we believe is a pragmatic emphasis on replicability. This paper presents two methods of estimating a study's replicability that researchers should consider reporting along with their statistical significant and effect size findings. One method of estimating the replicability of the findings deals with replication in the exact same system. The second method, which may contain subjective probability values, is used to estimate the replicability of a study's findings in a system that may differ from the initial system with respect to salient variables. The incorporation of the replicability estimates delineated in this paper would provide critical information to decision makers about the likelihood that the implementation of a particular method or treatment would produce similar results in their systems.

Déjà vu: Another Call for Replications of Research, Again

Most researchers would agree that the statistical significance levels and effects sizes reported in a study are important pieces of information (Fraas & Newman, 2000; Levin & Robinson, 2000; Robinson & Levin, 1997; Thompson, 1996, 1997, 1999a, 1999b). We take the position that what may be the most relevant piece of information practitioners and policy makers need to glean from a study is the ability of the study's findings to replicate (assuming the intent of the study is inferential not descriptive). If a study's findings are unlikely to replicate, the study's significance levels and effect sizes are virtually meaningless to interested practitioners and policy makers. Thus we as applied statisticians have the responsibility not only to provide estimates of a study's replicability but also delineate the assumptions on which these estimates are based.

In this paper we define two types of replications and present methods by which researchers can provide estimates of each type. The definition of replicability that we are developing in this paper is the extent to which a curriculum, treatment, etc. can be successfully implemented in two types of systems. One system is assumed to be an exact replication of the one used in the initial study. The replicability estimate that deals with this exact same system criterion is based on the same underlying assumptions of calculating the p value on a random sample of a known population. The other replicability estimate, assumes the system of interest is not an exact replication of the system used in the original study. This replicability estimate is based not only on random sampling assumptions but also on probability estimates, which will be somewhat subjective, that certain key variables will differ between the system used in the initial study and the system of interest. It is important to note that this second type of replication estimate can be calculated before as well as after the initial study is implemented. If it is calculated before the implementation of the study, it could assist in the re-design of the study before it is actually implemented. If calculated after the study is implemented, it will have implications for practitioners and decision makers.

Statistically Significant Exact Replications of a Study

One value we believe should be contained in research reports is the likelihood that the study's findings are replicable in the same system. Such a value should not take the place of statistical significance tests but rather should be reported along with them. We agree with Robinson and Levin (1997) who expressed the position that the probability value (p value)

produced by a statistical test is an important piece of information to report in a quantitative study. Robinson and Levin stated that "authors should *first* indicate whether the observed effect is a statistically improbable one (e.g., is the difference greater than what would be expected by change?)" (p. 22).

It is important, however, to not misinterpret a p value with respect to the likelihood of the replication of results (Nickerson, 2000). This point was addressed by Posavac (2002) who stated that "some believe [incorrectly] that rejecting a null hypothesis means that at least 95% of replications would be statistically significant (p. 102). Posavac does take the position, however, that rejecting a null hypothesis should increase the researcher's expectation that replications of the research would yield similar results.

Using p Values to Estimate Statistically Significant Exact Probabilities

Greenwald, Gonzalez, Harris, and Guthrie (1996) presented an analytic method by which a p value can be converted into a probability estimate that an exact replication of the research would produce a statistically significant result. Posavac (2002), who elaborated on the method proposed by Greenwald, et al., noted that an "*exact replication* means that the initial experiment is repeated using the same independent and dependent variables with the same number of participants selected in the same way from the same population" (p. 102). In this type of replication the difference between the replication and the original study is due to random variation. We believe that this is one type of replication that should be addressed by researchers.

Green et al. (1996) and Posavac (2002) suggest that the probability of a statistically significant exact replication (SSER) can be estimated from the probability of the statistical test. As a means of demonstrating how a researcher could convert a p value from a statistical test to an estimate of the SSER probability, we present a brief discuss of the procedure. It is beyond the scope of this paper to present the rationale on which this procedure is based. We encourage interested readers to review the works published by Greenwald et al. and Posavac for a more in-depth discussion of this concept.

An illustration

To illustrate the calculation of the SSER probability value, assume researchers are testing the difference between sample means of two independent groups consisting of 20 subjects each. Further assume that the t value produced by the difference between the two means recorded for their study was 2.150. Since this observed t value (t_{obs}) is greater than the two-tailed critical t

value (t_{crit}) of 2.024 for an alpha level of .05, the researchers would declare the difference between the two group means to be statistically significant. The question we believe is important for these researchers to address is: What is the chance that the difference between the two group means recorded for an exact replication of the study would be declared statistically significant?

Calculation of the SSER probability value. As noted by Posavac (2002), the probability of obtaining a SSER can be obtained by executing three steps. First, the replication t value (t_{rep}) is calculated by subtracting the critical t value used in the initial study from the study's observed t value. Thus the t_{rep} value is calculated as follows for our hypothetical example:

$$\begin{aligned} t_{rep} &= t_{obs} - t_{crit} \\ t_{rep} &= 2.150 - 2.025 \\ t_{rep} &= .125 \end{aligned}$$

Second, the researchers would obtain the one-tailed probability for this t_{rep} value of .125 with 38 degrees of freedom. With respect to the procedure used in this step, Posavac (2002, p. 108) stated that "a one-tailed test is used because one would want a replication to produce means of the same relative magnitudes as found in the first study." The one-tailed probability for the t_{rep} value of .125 with 38 degrees of freedom is .45.

Third, the researchers subtract the .45 probability value from 1.00, which produces a value of .55. This value indicates that the chance that an exact replication will be statistically significant is .55.

Points to Note Regarding the SSER probability value. Three points should be noted regarding this SSER probability value of .55. First, the SSER probability value is a function of the p value. However, practitioners need to be careful not to directly interpret the p value as a replicability value. Second, Greenwald et al. (1996) and Posavac (2002) recommended that SSER probability values should be considered upper limits. The reason for this recommendation is based on the fact that "even in a careful replication the participants would be a different sample from the population, the calendar date would be different, the weather would be different and so forth" (Posavac, p. 111).

Third, researchers may be surprised that for a study, such as the one used in our example, which had 38 degrees of freedom and an observed t value of 2.150 ($p = .038$), the chance that an exact replication will be statistically significant (SSER probability level = .55) is only slightly above the 50-50 level. In fact an observed t value for this hypothetical study would need to be

2.874, which produces a p value of .01, in order for the a SSER probability level to reach the .80 level.

To further emphasis this third point, a review of values produced by Posavac (2002) reveals that when degrees of freedom value is at least eight and the p value is .05, the SSER probability value will be .50. That is, there is a 50-50 chance of replicating significant findings. If the degrees of freedom value is at least eight and the p value is .01 for a two-tailed test, the SSER probability value will not be less than .73 or greater than .84. And if the degrees of freedom value is at least eight and the p value is .005, the SSER probability value will not be less than .80 and not greater than .92. (It is interest to note that these replicability values are less for corresponding p values for one-tailed tests.) Thus researchers need to be careful not to assume that statistically significant finds automatically mean that the chance of obtaining statistically significant exact replications for the study will be high. For this reason we believe that researchers should report the SSER probability value along with the probability of the observed t test.

Replication in a Different System

We believe that a second type of replication of findings is important for researchers to address. That is, the type of replication that deals with the question: Would the study's findings replicate in a different system than the one used in the initial study? It should be noted that we consider this type of replication of findings important even if an individual is interested in the same system in which the study was conducted, assuming the system is a dynamic one. That is, the system experiences considerable change with respect to the variables that may influence the replicability of the findings. Since most people attempt to relate research findings to systems that are different from their own or, at least, relate findings to systems that are similar but dynamic, we believe obtaining a likelihood estimate for this type of replication would be most valuable for them. The remaining portion of this section of the paper presents our *preliminary* attempt to develop a procedure for calculating such a likelihood estimate.

Estimate Procedure

The procedure we are proposing for the estimate of the likelihood of replication of findings for a different system than the one in which the study was conducted can best be presented through an example modeled on a study conducted by Benson, Aronson, Desmett, Shaheen, and Showalter (2002), which presented an evaluation of a multiage classroom

educational program. In our example we have children in grades 1-3 who were grouped in the same classroom and their teacher stayed with them for the three years. The evaluation indicated that the teachers volunteered for the project and were enthusiastic about the concept of multiage education. The project was supported fully by the principals and was enthusiastically supported by the parents. Achievement scores indicated moderate success of the project, as compared to national norms, and comparison students in the same school.

Internal validity issues are apparent, since parents voluntarily allowed their children into the project (see Campbell and Stanley, 1963, for a discussion of internal validity issues). Enthusiastic teachers might generate better results, no matter what the curriculum. In addition, a supportive principal might be partly (or entirely) responsible for the achievement results. Other internal validity concerns could also be raised.

External validity issues are also of concern with this study. Would the same effects be observed with less enthusiastic teachers? Will the same effects occur after the novelty of the multiage grouping wears off? Other external validity concerns could be raised (see Campbell and Stanley, 1963, for a discussion of external validity issues). The concept of replicability, though, is different than internal validity and external validity. It is based on the realization that any implementation is accomplished in a system, and the realization that that system is likely to be dynamic. We believe that the likelihood of replicating a study's findings in a different system or even the same *dynamic* system is crucial to estimate.

Important variables. The first step in the estimation process is to identify key variables that influenced the findings but maybe different in the new system. Let us assume that for our multiage project example four such variables were identified:

1. Twelve volunteer teachers were used.
2. The study involved supportive principals.
3. The study used 240 volunteer (supportive) parents.
4. A total of 5 days of in-service training was given to the teachers on the multiage project.

As an illustration of how these variables could influence the replicability of the findings of the original study, consider the principals variable. If a principal leaves, the project will, in all likelihood, be supported less by the new principal. The new principal may even kill the project, not because the project is ineffective, not because the concept of multiage education is bad, but

because the crucial component of the system (the principal) does not believe in or want the project.

The likelihood of each crucial component changing should be taken into account when the project is envisioned. If a particular component is likely to change, then the project should be devised so it is immune to that change--in the case of principal change--the project should be made "principal proof."

Once the variables are identified, the second step is to estimate the proportion of the R^2 value accounted for by each variable, the probability of that variable changing, and the probability of the changed variable being negatively influential on the original findings. Table 1 contains such hypothetical values of these estimates for our example.

 Insert Table 1 about here

The proportion of the R^2 value accounted for by each component is determined. This could be accomplished with GLM if enough implementation sites were available (similar to meta analysis), or conceptualized either before the study started or afterwards. In the example, here we provide "educational guesses." For instance, it is likely that some teachers will leave the project. Some may become disillusioned with the project or with education in general. Others may find a more lucrative job in another district or another profession. Nevertheless, other enthusiastic teachers are likely available, so the systemic effect on the project of teacher change would be minimal.

On-the-other hand, the likelihood of a principal leaving the system is high (estimated to be .70 in a three-year period) and the likelihood of the replacement being equally enthusiastic is low (.40). Indeed, most replacement principals may gut the project, leading to absolutely no replicability from the component of the principal. Therefore, because of the high probability of principal change, and high probability of a different (lower) level of support, the overall replicability is lowered.

Parent turnover will be at least 33% every year, with third graders moving to fourth grade. But we suspect that the parents of the incoming first graders will be just as enthusiastic (maybe even more so if the project is a success). Thus, the high turnover rate (large system change) of parents will have little effect on replicability--the project is "parent proof." If the

staff development is “packaged” then it could easily stay the same from year to year. This part of the system would likely be stable.

Actual probabilities may be quite difficult to determine. To deal with this problem, one might rate the stability of each component on a 1 to 5 scale, with 5 being the most stable. Such estimates and the calculation of the reliability value for the multiage example are listed in Table 2. It should be noted that a replicability value calculated in this manner would produce higher values the more stable key variables are from the system used in the initial study and the system of interest especially for the variables that account for the higher proportion of the R^2 value.

 Insert Table 2 about here

Estimating replicability before implementing. If a researcher calculated replicability before first implementing a new project, and obtained a low replicability value, the researcher might try to re-conceptualize the project by either doing something to minimize system change or to minimize the effects of the change in any one component. One could minimize system change in the multiage project by getting the school board to mandate multiage in all elementary schools, or find another district where all the principals are supportive of multiage programs. It should be noted that a replicability value calculated in this manner would produce higher values the more stable key variables are from the system used in the initial study and the system of interest especially for the variables that account for the higher proportion of the R^2 value.

Minimizing the effects of change in teachers could be accomplished by each project teacher identifying a non project teacher who would like to be in the project and then keeping that teacher informed about multiage groping during the year. This “information partnership” actually becomes a new component of the project (or at least modifying the teacher component.) Curricula that purport to be “teacher proof” are another example of minimizing the effects of teacher change.

If the replicability index is low and the researcher cannot identify changes or strategies that would make it higher then the project should not be implemented. The time of teachers, principals, parents, staff developer, and especially students should not be wasted. If there is very little hope for replication of a particular project, then we have no business investigating the effectiveness of that project.

How about changes in system components not relevant to the project? Changes in components that are not relevant to the project will not affect the replicability, by definition. Nor will these changes affect the index, as the percent of variance accounted for is 0 and the contribution of that component would be 0. Unfortunately, in most educational systems, many components can influence the success of a project.

Implications

The implications of this paper relate our position that statistical significance and effect size are important concepts but they must be examined in light of replicability. Replicability is and of itself not a one dimensional concept but rather a multi-dimensional one. In this paper we identified two types of replication estimates. The first type is the SSER probability estimate, which is based on traditional statistically assumptions and probability concepts. The second type is related to design and subjective probability issues. This approach provides a number of advantages. First, it can assist in the teaching of research design. That is, teaching this replication estimate emphasizes the need for researchers to attempt to identify the relevant variables in a study. Second, it can improve communication among researchers regarding relevant variables in a study in order to improve the design of such studies. Third, it encourages the use of meta-analysis to identify relevant variables. Fourth, it provides a method of simulating the effects the relevant variables on replicability of findings. One can simulate small changes or large changes on relevant variables and the impact of these changes on replicability. As one can see, this second estimate is not a static approach but rather a dynamic one and may only be limited by the investigated creativity and in sight.

With respect to methodology, more partial replications should be encouraged. Partial replication can be conducted by one of two approaches. First, half of the study could be an exact replication, and the other half could be an extension (into another grade level, using different in-service materials, or checking on efficacy in another bureaucratic situation). Second, the researcher could put a slight twist on the implementation, by reducing or eliminating a component, shortening the period, streamlining in-service, or monitoring more closely the actual implementation.

With respect to analysis, weighting coefficient from a previous study can be applied to the current results. Alternatively, the weights from the analyzed sample can be cross-validated

with the holdout sample. At the least, the effect size should be compared with the effect sizes found in the knowledge base.

With respect to discussion of how a study's results fit into the knowledge base, the researcher should do just that. Too often, the discussion section (at least in dissertations) is written in a compressed period, and that section is given short shrift. The researcher should admit here the shortcomings of the study, and identify what "camp" the researcher is a part of.

References

- Benson, S. N. K., Aronson, E., Desmett, P., Shaheen, M., & Showalter, J. (2002, October). *Multiage education: A process and product evaluation*. Paper presented at the annual meeting of the Midwestern Educational Research Association, Columbus, OH.
- Campbell, D. T. & Stanley, J. C. (1963). *Experimental and quasi-experimental designs for research*. Chicago, IL: rand McNally College Publishing Co.
- Fraas, J. W., & Newman, I. (2000, October). *Testing for statistical and practical significance: A suggested technique using a randomization test*. Paper presented at the annual meeting of the Mid-Western Educational Research Association, Chicago, IL.
- Greenwald, A. G., Gonzalez, R., Harris, R.J., & Guthrie, D. (1996). Effect sizes and p values: What should be reported and what should be replicated? *Psychophysiology*, *33*, 175-183.
- Levin, J. R. & Robinson, D. H. (2000). Rejoinder: Statistical hypothesis testing, effect-size estimate, and the conclusion coherence of primary research studies. *Educational Researcher*, *29* (1), 34-36.
- Nickerson, R. S. (2000). Null hypothesis significance testing: A review of an old and continuing controversy. *Psychological Methods*, *5*, 241-301.
- Posavac, E. J. (2002). Using p values to estimate the probability of a statistically significant Replications. *Understanding Statistics*, *1*(2), 101-112.
- Robinson, D. H. & Levin, J. R. (1997). Reflections on statistical and substantive significance, with a slice of replication. *Educational Researcher*, *26*(5), 21-26.

Thompson, B. (1996). AERA editorial policies regarding statistical significance testing: Three suggested reforms. *Educational Researcher*, 25 (2), 26-30.

Thompson, B. (1997). Editorial policies regarding statistical significance tests: Further comments. *Educational Researcher*, 26 (5), 29-32.

Thompson, B. (1999a). Statistical significance tests, effect size reporting and the vain pursuit of pseudo-objectivity. *Theory and Psychology*, 9 (2), 191-196.

Thompson, B. (1999b). Journal editorial policies regarding statistical significance tests: Heat is to fire as p is to importance. *Educational Psychology Review*, 11 (2), 157-169.

Table 1

Proportion of the R^2 Accounted for in the Dependent Variable
and the Probability Values for Each Variable

Variable	Proportion of R^2	Estimate of the Probability of Change	Estimate of Probability of Negative Impact
Teacher	.50	.30	.03
Principals	.20	.70	.60
Parents	.10	.33	.02
Staff Development	.20	.02	.02

Table 2

Calculation of the Replicability Value			
Variable	Proportion of R^2	Stability	(Proportion of R^2) * (Stability)
Teacher	.50	4	$.50 * 4 = 2.00$
Principals	.20	1	$.20 * 1 = 0.20$
Parents	.10	5	$.10 * 5 = 0.50$
Staff development	.20	5	$.20 * 5 = \underline{1.00}$
			Replicability = $3.70 / 5 = .74$



U.S. Department of Education
Office of Educational Research and Improvement (OERI)
National Library of Education (NLE)
Educational Resources Information Center (ERIC)



REPRODUCTION RELEASE

(Specific Document)

TM034969

I. DOCUMENT IDENTIFICATION:

Title: Deja vu: Another Call for Replications of Research, Again	
Author(s): Isadore Newman, Keith McNeil, and John W. Fraas	
Corporate Source: The University of Akron	Publication Date: April 21, 2003

II. REPRODUCTION RELEASE:

In order to disseminate as widely as possible timely and significant materials of interest to the educational community, documents announced in the monthly abstract journal of the ERIC system, *Resources in Education* (RIE), are usually made available to users in microfiche, reproduced paper copy, and electronic media, and sold through the ERIC Document Reproduction Service (EDRS). Credit is given to the source of each document, and, if reproduction release is granted, one of the following notices is affixed to the document.

If permission is granted to reproduce and disseminate the identified document, please CHECK ONE of the following three options and sign at the bottom of the page.

The sample sticker shown below will be affixed to all Level 1 documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL HAS BEEN GRANTED BY

_____ Sample _____

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

1

Level 1



X

Check here for Level 1 release, permitting reproduction and dissemination in microfiche or other ERIC archival media (e.g., electronic) and paper copy.

The sample sticker shown below will be affixed to all Level 2A documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE, AND IN ELECTRONIC MEDIA FOR ERIC COLLECTION SUBSCRIBERS ONLY, HAS BEEN GRANTED BY

_____ Sample _____

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

2A

Level 2A



Check here for Level 2A release, permitting reproduction and dissemination in microfiche and in electronic media for ERIC archival collection subscribers only

The sample sticker shown below will be affixed to all Level 2B documents

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE ONLY HAS BEEN GRANTED BY

_____ Sample _____

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

2B

Level 2B



Check here for Level 2B release, permitting reproduction and dissemination in microfiche only

Documents will be processed as indicated provided reproduction quality permits. If permission to reproduce is granted, but no box is checked, documents will be processed at Level 1.

I hereby grant to the Educational Resources Information Center (ERIC) nonexclusive permission to reproduce and disseminate this document as indicated above. Reproduction from the ERIC microfiche or electronic media by persons other than ERIC employees and its system contractors requires permission from the copyright holder. Exception is made for non-profit reproduction by libraries and other service agencies to satisfy information needs of educators in response to discrete inquiries.

Sign here, →
ease

Signature:	Printed Name/Position/Title: John W. Fraas, Ph.D. Trustees' Professor		
Organization/Address: Ashland University, 401 College Avenue, Ashland, Ohio, 44805	Telephone: 419-289-5930	FAX: 419-289-5910	Date: 5/20/2003
	E-Mail Address: jfraas@ashland.edu		



III. DOCUMENT AVAILABILITY INFORMATION (FROM NON-ERIC SOURCE):

If permission to reproduce is not granted to ERIC, or, if you wish ERIC to cite the availability of the document from another source, please provide the following information regarding the availability of the document. (ERIC will not announce a document unless it is publicly available, and a dependable source can be specified. Contributors should also be aware that ERIC selection criteria are significantly more stringent for documents that cannot be made available through EDRS.)

Publisher/Distributor:
Address:
Price:

IV. REFERRAL OF ERIC TO COPYRIGHT/REPRODUCTION RIGHTS HOLDER:

If the right to grant this reproduction release is held by someone other than the addressee, please provide the appropriate name and address:

Name:
Address:

V. WHERE TO SEND THIS FORM:

Send this form to the following ERIC Clearinghouse: University of Maryland ERIC Clearinghouse on Assessment and Evaluation 1129 Shriver Laboratory College Park, MD 20742 Attn: Acquisitions
--

However, if solicited by the ERIC Facility, or if making an unsolicited contribution to ERIC, return this form (and the document being contributed) to:

**ERIC Processing and Reference Facility
1100 West Street, 2nd Floor
Laurel, Maryland 20707-3598**

Telephone: 301-497-4080

Toll Free: 800-799-3742

FAX: 301-953-0263

e-mail: ericfac@inet.ed.gov

WWW: <http://ericfac.piccard.csc.com>