

DOCUMENT RESUME

ED 041 943

TM 000 035

AUTHOR Lewy, ArieH
TITLE Item Analysis Based on Classroom Performance.
PUB DATE Mar 70
NOTE 9p.; Paper presented at the annual meeting of the American Educational Research Association, Minneapolis, Minn., March 1970

EDRS PRICE EDRS Price MF-\$0.25 HC-\$0.55
DESCRIPTORS *Achievement Tests, Correlation, *Group Norms, Group Tests, *Item Analysis, *Statistical Analysis, Statistics, Tests

ABSTRACT

Statistical procedures employed in item analysis are based on the performance of the individual on a given test. If, however, one desires to assess the efficiency of some treatment, the researcher is interested in the performance of groups of students who received the treatment rather than in the performance of any one individual. In such cases statistics based on individual performance do not constitute a proper basis for item selection and for interpretation of test results. In this context, therefore, the discrimination power of an item should be conceived as its capacity to discriminate among classes and not among individuals. This approach is not only logically sound, but also increases the quantity and quality of the information available. Three new parameters are introduced: item-test correlation based on class data, item standard deviation based on class data, and the intraclass correlation coefficient of an item. An item analysis of a fourth grade arithmetic test is used as an example. Other class data parameters are suggested for future research. (DG)

ITEM ANALYSIS BASED ON CLASSROOM PERFORMANCE

Dr. ARIEH LEWY

Tel-Aviv University

Department of Education

U.S. DEPARTMENT OF HEALTH, EDUCATION
& WELFARE
OFFICE OF EDUCATION
THIS DOCUMENT HAS BEEN REPRODUCED
EXACTLY AS RECEIVED FROM THE PERSON OR
ORGANIZATION ORIGINATING IT. POINTS OF
VIEW OR OPINIONS STATED DO NOT NECES-
SARILY REPRESENT OFFICIAL OFFICE OF EDU-
CATION POSITION OR POLICY.

The Background

Statistical procedures employed in test analysis and especially in item analysis are based on the performance of the individual in a given test. Item characteristics which are of interest for the evaluation expert and which frequently serve as a basis for item selection are: the difficulty level of an item, its correlation with the scores on the total test, its correlation with other items of the same test and with one or with several external criteria scores. Several standard computer programs ⁽⁵⁾ report also the standard deviations of each item, which in most multiple choice tests are curvilinear function of the item difficulty and thus do not convey any new information which is not present in the item difficulty index itself. That test statistics are based on measurements of individual performance is possible due to the fact that the roots of statistical test theory are in clinical testing, which, in turn, has its major interest in individual differences. Whenever one wishes to develop an instrument which will maximize the discrimination among individuals with regard to certain trait or capability such procedures may be most efficient. But the concern of educational testing and measurement is not only the assessment of individual performance. Tests are used within the framework of research studies come to assess the efficiency of some treatment, like teaching method, curriculum material, organizational setting etc. employed in classroom situations to groups of students. In such cases the researcher is interested in the performance of groups of students who underwent some treatment rather than in the performance of any individual. In such cases statistics based on individual performance do not constitute a proper basis for item selection and for interpretation of test results. Whenever tests

Read at the Annual Meeting of the American Educational Research Association,
March 5, 1970, Minneapolis

ED041943

TM 000 035



are used for the purpose of assessment of some aspects of classroom teaching item analysis has to be performed on the basis of differences between achievements of classes.

This approach can be justified on the basis of logical considerations but additionally it possesses empirical advantages by providing a variety of information which is not available whenever individual data are used for item analysis.

Logical basis for using class-performance data

The fact that in classroom experimentation the unit of observation is not the individual but the class has been pointed out by Lindquist's classical book on Experimental Design⁽⁶⁾. Lindquist calls attention to the fact that performance of students within the same class are not independent one of the other, since they are affected by circumstances prevailing in the class. Therefore the number of independent replications is equal to the number of classes included in the sample and not to the number of individuals comprised in the sample. This idea has been further developed by Lord⁽⁷⁾ and utilized in several educational studies like Bock⁽¹⁾, but surprisingly has never been applied to item analysis and item selection procedures. Since the natural setting of test administration is the classroom one has to consider the whole class as a unit and compare performances of classes rather performances of individual. Accordingly questions usually asked in the context of item analysis have to be reformulated. With this view in mind the discrimination power of item e.g. has to be conceived as its capability to discriminate among classes and not among individuals. For this reason it is proposed that the discrimination power of an item should be defined as the correlation between the average achievement of students in each class on a particular item with the average performance of the class on the whole test. Thus far the logical reasons for the suggested modification of item analysis have been described and a new definition for the concept "item discrimination" power has been suggested. But this approach not only increases the logical precision of a concept widely used, but also increases the amount and the quality of information which is processed by item analysis.

Empirical Advantages

The increase in amount and the improvement of the quality of information is illustrated by the presentation of results obtained from item analysis of a fourth grade arithmetic test.

The Test

A 42 item test has been constructed to measure achievement of fourth grade pupils in arithmetic⁽⁴⁾. The test covered several topics taught in the 4th grade according to the local syllabus. In addition to topics taught in all schools the test contained two subsets of items which are of special interest for the purpose of this paper. The first of them is Fractions. This topic has been taught in grade four of some of the schools which study arithmetic according to a new experimental program, while in other schools it is taught only in grade five. Thus some classes studied the topic Fractions before taking the test while others did not systematically study this topic. The second subset of items labelled in this paper, Theory, deals with general aspects of mathematics and contains questions like "Each number which can be divided by 4, can be divided by 2 too" etc., a topic which has not been dealt with directly in classes and may be considered as a kind of measure of numerical ability. And finally there is a third set which contains a single item - item 12 - labelled error which in this case refers to the fact that the correct response has been erroneously provided for the computer analysis and thus the results of this item have to be disregarded, but for analytical purposes it turned out to be valuable for the present topic. The test has been administered to 107 classes comprising 3057 students with an average of 28 students per class. The classes represented the whole range of differences in ability levels existing in the country.

The results of item analysis

Some results of the analysis are presented in table 1. This table contains information based traditionally on individual data, and information based on class data. It can be seen that the test is a typical one, its items represent different difficulty levels.

Their distribution is as follows:

<u>Percent of students giving correct response</u>	<u>No. of item</u>
20 - 29	4
30 - 39	2
40 - 49	9
50 - 59	7
60 - 69	8
70 - 79	7
80 - 89	8

On most of the items 40 - 80 percent of correct responses were obtained, but some items are rather difficult and others quite easy.

Item-test correlation

The traditionally reported item total test biserial correlation coefficient is listed in column (2) of table 1. In column 4 the correlation coefficients among the same variables based on class performance data is listed. The first striking observation is that correlations based on class performance are considerably higher than those based on individual performance data. However a constant difference between the two types of correlations would not be of interest. But one can see that the range of correlations based on class performance is higher than of those based on individual data. The range of the point biserial correlations is from .2 to .6 with a median value of .4 while the range of class-correlations occupies almost the whole possible range of positive correlations, it runs from .2 to .9 with a median value of .7. The reason for this difference is that individual performances are scored dichotomously (0 or 1) while class averages for each item can carry a variety of fraction values from 0 to 1.

It should be added that in the present table the individual point-bis correlations are computed without correction formula for auto correlation and thus their values are slightly inflated⁽⁸⁾ while the class correlations are computed with regard to the sum of other items only. Thus the difference between the actual correlation coefficients is larger than those presented here.

It can be seen that in using class data we obtain a coefficient which has more refined discriminational value than that obtained by using individual data.

Standard deviation of class averages

Using class data also enables us to examine some characteristics of items which are not described at all in conventional item analysis.

I refer here to the frequency distribution of class-averages which is presented on the right part of table 1. It can be seen e. g. that on item no. 1 two classes obtained 50-59 percent correct responses and 13 classes obtained 60-69 percent responses. On the same item 21 classes obtained 100 percent of correct responses, i. e. in 21 classes all students responded correctly to this item. The distribution of class averages is characterized by a single parameter i. e. the standard deviation of class averages. Unlike the standard deviation based on individual data, the standard deviation of class averages (see col. 3) is not a function of the item difficulty and thus it conveys information not contained in the item difficulty index itself. Theoretically it is possible that one item having a given difficulty level has widely different averages in a variety of classes while another item with the same difficulty level has quite similar averages in most classes. Thus e. g. in the data presented here both items 37 and items 41 have a difficulty level of 45 and a slightly different standard deviation.

The differences in the size of standard deviation in the present sample are moderate. However, considerable differences have been detected by the author in other tests analyses.

Interclass correlation

A third type of information which may be of interest for the evaluation expert is the intraclass correlation coefficient reported in column 6 of the handout (Haggard, 1958)⁽³⁾. This parameter is the function of the ratio obtained by the division of the between classes mean sum of square by the within classes mean square. The intraclass correlation is an index of the degree of similarity of responses given by pupils in any single class. As indicated already the topic Fraction has been taught in some classes and in other classes it has not been taught at all. As a result of this the intraclass correlations of items in this category are relatively high. Thus e.g. the intraclass correlation of item 41 is very high. This item represents a relatively easy question in fraction. Those classes who studied the relevant material obtained very high averages on this item while classes which did not study this material obtained averages close to zero. The frequency distribution of class means on this item has two peaks i. e. it has bimodal shape. This is an interesting example of an item whose variance stems mainly from between class differences while on the other hand the variance of item 33 stems mainly from within class differences. The former one is a content oriented item while the later one is more saturated with general ability. The comparison of these two items supports the feeling that items testing achievement have higher intraclass correlation than items measuring ability.

Space does not permit to present other statistics which can be computed on the basis of class data. But the purpose of the present paper is not to convey the findings of a specific study but rather to call attention to the existence of basic data which have been completely neglected in dealing with item statistics and which may yield a wide variety of parameters increasing the amount and the better quality of information obtained from test results. It is often thought that the currently based statistics tell everything about test items, which may be of interest. In fact the practice of item analysis and item selection has not changed over decades and the endeavor of psychometrists is more concentrated on provisions of computational shortcuts for the existing parameters than for seeking additional parameters which may increase our understanding of tests.

Topics for Further Study

The present paper introduces three parameters which have not been used before and these are: 1) item-test correlation based on class data; 2) item standard deviation based on class data, and 3) intraclass correlation coefficient of the item. In addition to these parameters some other topics for research utilising class data are as follows:

1. Inter-item correlation matrix computed on basis of class results. It is expected that these correlation coefficients will have a pattern different than those based on individual data.
2. Factor analytic studies can be performed on both types of correlation matrices.
3. Test-reliability coefficients of various types can be computed on basis of class results, too. Thus it is possible to employ split-half computational procedures as well as the formulae of Kude Richardson generalized by Cronbach⁽²⁾ as the coefficient alpha.
4. Studies can be conducted to identify factors which affect the similarity or the difference between statistics obtained from individual and from class data.
5. Studies can be conducted to find out whether items testing achievement can be differentiated from items testing ability by the size of their intraclass correlation.

It can be envisaged a long series of studies based on the idea mentioned here, which may change the nature of major topics dealt with in test theory. Studies based on class-data may contribute new chapters to standard text books on measurement and evaluation and may sharpen the understanding of differences between tests used for the purpose of individual guidance and selection on one hand, and tests used for the purpose of testing the efficiency of classroom teaching practices on the other hand.

BIBLIOGRAPHY

1. Bock, R. D. "Contributions of Multivariate Experimental Design to Educational Research" In Cattell B. R. (ed) Handbook of Multivariate Experimental Psychology. Rand M. Nally, Chicago 1966.
2. Cronbach, L. J. "Coefficient alpha and the internal structure of tests". Psychometrika 1951, v. 16 p. 297-334.
3. Haggard, E. A. Intraclass Correlation and the analysis of variance New York; Dryden 1958
4. Lewy A. and Chen M. "Scholastic Achievements of Fourth Grade Students in Israel" (in preparation)
5. Lewy A. and Crawford R. "Scoring Test Battery" Educational and Psychological Measurement. v. 26 p. 185-188, 1966
6. Lindquist E. F. Design and Analysis of Experiments in Psychology and Education. Boston: Houghton Mifflin 1963.
7. Lord F. M. "Test Norms and Sampling Theory", Jour. of Exp. Ed., v. 27, p. 247-263.
8. Wolf, R. (To be supplied)
Journal of Educational Measurement, 1967 (?)