

HOW WE SHOULD MEASURE "CHANGE"—OR SHOULD WE?

LEE J. CRONBACH¹ AND LITA FURBY²

Stanford University

Procedures previously recommended by various authors for the estimation of "change" scores, "residual" or "basefree" measures of change, and other kinds of difference scores are examined. A procedure proposed by Lord is extended to obtain more precise estimates, and an alternative to the Tucker-Damarin-Messick procedure is offered. A consideration of the purposes for which change measures have been sought in the past leads to a series of recommended procedures which solve research and personnel-decision problems without estimation of change scores for individuals.

A persistent puzzle in psychometrics has been "the measurement of change." Many investigators have felt, for reasons good or bad, that their substantive questions required a measure of gain in ability or shift in attitude. "Raw change" or "raw gain" scores formed by subtracting pretest scores from posttest scores tend to fallacious conclusions, primarily because such scores are systematically related to any random error of measurement. Although the unsuitability of such scores has long been discussed, they are still employed, even by some otherwise sophisticated investigators.

At the end of this paper the authors argue that gain scores are rarely useful, no matter how they may be adjusted or refined. The authors also distinguish four kinds of inquiry for which such scores have been used, and conclude that only one of these purposes is well served by any kind of gain score. This argument applies not only to changes over time, but also to other differences between two variables.

The first part of the paper proposes superior ways of estimating true change and true residual change scores. It may seem pointless to discuss such matters when in the end we recommend against their use (save in a few kinds of investigations). However, the development of formulas clarifies the model, providing a

¹ Initial work on this paper was conducted under a contract with the Cooperative Research Program of the United States Office of Education. We thank many colleagues for helpful comments, particularly Janet D. Flanshoff, Gene V. Glass, Ingram Olkbi, and Julian C. Stanley, Jr.

² Requests for reprints should be sent to Lee J. Cronbach, School of Education, Stanford University, Stanford, California 94305.

³ Now at the University of Poitiers, France.

ing the person's "true" status at these times, are postulated. The "true difference" D_{∞} equals $Y_{\infty} - X_{\infty}$. The key topic of the McNemar paper and the Lord papers is the determination of regression coefficients for an estimator of the form:

$$\hat{D}_{\infty} = \beta_1 X + \beta_2 Y + \text{constant} \quad [1]$$

As Lord has written more extensively than McNemar on this matter, we shall refer only to Lord hereafter, though many of the comments also apply to what McNemar has said.

The development holds so long as X and Y are referred to the same numerical scale. That is to say, the investigator must be willing to say what score on Y (or Y_{∞}) is comparable to a given score on X (or X_{∞}). The relation must be reciprocal in the sense that, if the given X is mapped into a certain Y , that value of Y is mapped into the given X . It is particularly important to note that this formulation sidesteps the philosophically troublesome question, "Are pretest and posttest 'measuring the same variable'?" A common metric is the only requirement. Even this requirement is dispensed with when we turn to a regression estimate of outcome.

The model applies to any two measures that can be sensibly expressed on the same numerical scale. This might be, for example, a standard-score scale or an age-equivalence scale. The data might be ratings of two distinct traits expressed on the same reference scale. Hence statements made about change scores can be extended to any kind of difference score. Among the more famous lines of research employing difference scores are work on "overachievement," "empathy and insight," "self-concept," versus "ideal self," and "differential aptitudes." Indices of the same psychometric character are also involved in studies of retention and transfer.

Assumptions. There are two variables and X Y . For each variable, the expected value over independent observations of the same person defines a true score: $E(X) = X_{\infty}$ and $E(Y) = Y_{\infty}$. Then $D_{\infty} = Y_{\infty} - X_{\infty}$. The investigator who identifies Y with X as "the same operation" must keep the true scores distinct. In a study of change, X_{∞} is the person's average score over measurements of X that might be made at Time 1; Y_{∞} is the average over observations by the same procedure at Time 2.

Errors are uncorrelated with true scores. But we do not assume that all correlations are equal. Sometimes X and Y observations are "linked," as when the two scores are obtained from a single test or battery administered at one sitting, or when observations on different occasions are made by the same observer. The correlation between linked observations will ordinarily be higher than that between independent observations. This distinction has not been made in the literature on change, but it does appear in a paper by Stanley (1967) on difference scores.

We develop a formal mathematical model so that some parts of our argument can be explicit rather than intuitive. We retain the classical concept of strictly parallel observations, but modify the classical concept of independence, as suggested above.

1. Let X_{it} represent the observed score of person i on the X variable observed under condition t . The condition t may be, for example, a particular form of the test that measures X . In studies where there are several sources of "error" such as test form, observer, and short-term fluctuations in the state of the person, we assume that these sources are completely confounded in the design used to determine reliability coefficients and correlations.

2. Where the classical theory considers true score as the sum of observed score and error, we introduce two random errors: $X_{it} = X_{i\infty} + e_{it} + f_{it}$. Likewise, $Y_{it} = Y_{i\infty} + e_{it} + f_{it}$. This is required to formalize the concept of independence adequately.

If, as in the classical concept of parallel measures, one assumes that the several measures of X (or Y) are strictly interchangeable, there can be no linkage. Machinery for explicitly defining independence can be developed with some concepts from generalizability theory. There is a universe of possible conditions of observation of X . (While observations may vary with respect to test form, occasion, observer, etc., we shall avoid here the complications of multifacet theory [Gleser, Cronbach, & Rajaratnam, 1965].) That theory recognizes, as this paper does not, that one may have two measurements with the same test form on different occasions, or two measurements on different test forms on the same occasion, or

two measurements where both form and occasion differ.) There is a universe of observations of Y , and we shall assume that these may be made under the same set of conditions f that are used for X observations. Then observations X_i and Y_i made under the same condition i are said to be linked, and observations X_i and Y_j made under different conditions are said to be independent.

Formally, the model specifies that conditions of observations are drawn from the universe. If a single i is drawn and used to obtain both scores $X_{i\beta}$ and $Y_{i\beta}$ for person β , we have linkage; if i and j are drawn independently to give scores $X_{i\beta}$ and $Y_{j\beta}$, we have independence. Errors are random and independent, except that when X and Y are both observed under condition i , the f and i components are sampled simultaneously. It will not be true, in general, that $\sigma(f_{i\beta}, f_{j\beta}) = 0$. The following assumptions are made regarding all e components: Their mean over persons is zero for every condition; their variance over persons is the same for every condition; their intercorrelations with other components are zero. The e components satisfy the same assumptions and $\sigma(e_{i\beta}, e_{j\beta}) = \sigma(f_{i\beta}, f_{j\beta}) = 0$. With regard to the f components, zero means, equal variances, and zero correlations with true scores are assumed (likewise for f). Also, $\sigma(f_{i\beta}, f_{j\beta}) = 0$ for all pairs where i is not identical to j .

3. The measures of X made under different conditions are parallel. It follows from the assumptions above that the measures have equal means, equal variances, and equal intercorrelations. The same is true for Y measures.

4. It follows, now, that

$$\sigma(X_{i\beta}, Y_{j\beta}) = \sigma(X_{i\beta}, Y_{i\beta}),$$

hence this covariance is the same for all independent observations of X and Y .

For linked observations the covariance

$$\sigma(X_{i\beta}, Y_{i\beta}) = \sigma(X_{i\beta}, Y_{i\beta}) + \sigma(f_{i\beta}, f_{i\beta})$$

We assume that $\sigma(f_{i\beta}, f_{i\beta})$ is the same for all i , hence that the linked covariance is the same for all linked X, Y observations. The covariance of f with f may be large or small, depending on the extent to which the condition i influences the X and Y performances.

We shall simplify and compress notation in several convenient ways. We write ρ_{XY} for a reliability coefficient, and ρ_{XY} for a correlation.

For emphasis, when linked observations are correlated or used to form a difference score we may refer to X_i and Y_i ; their covariance and correlation we shall designate ρ_{XY} and ρ_{XY} . We shall similarly write X_i and Y_i (or the like) for a pair of independently observed scores, and for the covariance and correlation shall write ρ_{XY} and ρ_{XY} .

Ordinarily, $|\rho_{XY}| > |\rho_{XY}|$. It is possible that $\rho_{XY} < \rho_{XY}$, where X_i and Y_i have some complementary relation. For instance, if rate of reading and comprehension are measured on the same selections simultaneously, one score is likely to rise at the expense of the other and $\rho(f_{i\beta}, f_{j\beta})$ will be negative.

5. It is assumed that the population parameters are known. All parameters considered must be for the same population.

A regression estimate of a true score is usually improved if the data permit one to derive an equation for a subpopulation—for example, for ninth-grade boys in a certain school rather than for the national ninth-grade population. The investigator using actual data will often have made no reliability study on his sample. If he uses a reliability coefficient or a value of ρ_{XY} from the published study he must adjust it to take into account the variance of his sample.

RELIABILITY OF A DIFFERENCE SCORE

As Stanley (1967) pointed out, linkage must be taken into account in defining a reliability coefficient for differences. Classical theory, ignoring linkage, defines a reliability

$$\rho_{DD} = \frac{\sigma^2_{D_{obs}} + \sigma^2_{\rho_{XY}} - 2\sigma_{\rho_{XY}}}{\sigma^2_X + \sigma^2_Y - 2\sigma_{XY}}, \quad [2]$$

where $D = Y - X$. The reliability coefficients ρ_{XY} and ρ_{XY} are correlations of independent observations. Likewise, because of the independence assumption, classical theory can only interpret ρ_{XY} as what we have called ρ_{XY} .

The reliability of a difference score is defined as the correlation of the score with an independently observed difference. Unlike the classical papers, Stanley distinguishes $\rho(\sigma^2_{X_{i\beta}}, \sigma^2_{Y_{i\beta}})$ from $\rho(\rho_{XY}, \rho_{XY})$. These are distinct kinds of reliability, one for a difference of linked X and Y and one for a difference of independent X and Y .

If the observed difference is $D_{obs} = Y_i - X_i$

Correction by simple regression for error in X . If \hat{X}_i is a regression estimate of X_i from X ,

$$D_{obs} = Y_i - \hat{X}_i = Y_i - \rho_{XY}X_i + \text{constant} \quad [9] [2]$$

In this and all other equations through Equation 20, the constant is one-Sigma regressed the mean (over persons) of the estimates equal to the mean raw gain. The foregoing estimate has occasionally been suggested (e.g., Trimble & Cronbach, 1943), but it has seen little use. Closely related concepts appear in Lord's comments on analysis of covariance (1960) and in a series of Swedish papers on the effect of schooling on intelligence (see Harnqvist, 1968).

Correction by simple regression for error in X and in Y . To take a further step, let \hat{Y}_i be an estimate of Y_i from Y . Then

$$D_{obs} = \hat{Y}_i - \hat{X}_i = \rho_{YY}Y_i - \rho_{XY}X_i + \text{constant} \quad [10] [3]$$

This procedure does not take the X, Y correlation into account.

The Lord Procedure

Lord pointed out that unless $\rho_{X_i} = 0$, both X and Y yield information about $X_{i\beta}$. A multiple regression procedure can be used to obtain

$$\hat{X}_{i\beta} = \rho_{XX}X_i + \beta(X_{i\beta}, Y_i)(Y_i - X_i) + (1 - \rho_{XX})\bar{X} \quad [11]$$

Here $Y - X$ is a partial variate, the deviation of Y from the value predicted by the regression of Y on X in the population to which the other parameters in Equation 11 apply.

For the linked case we have

$$Y_i - X_i = Y_i - \frac{\rho_{XY}}{\sigma^2_X}(X_i - \bar{X}) - \bar{Y} \quad [12]$$

We know that

$$\beta(X_{i\beta}, Y_i - X_i) = \frac{\sigma_{X_i(Y_i - X_i)}}{\sigma^2(Y_i - X_i)} \quad [13]$$

Now $\sigma_{X_i} = \sigma_{X_i} = \rho_{XY} = \rho_{XY}$ (not ρ_{XY}). Using Equations 12 and 13, the numerator of β becomes

$$\sigma_{X_i(Y_i - X_i)} = \sigma^2(Y_i - X_i) \rho_{XY} \quad [14]$$

and the denominator becomes

$$\sigma^2(Y_i - X_i) = \sigma^2_Y(1 - \rho^2_{XY}) \quad [15]$$

Substituting in the regression equation, we

(i.e., if X and Y are experimentally independent), the covariance with an independently observed difference D_{obs} is

$$\sigma^2_{\rho_{XY}} + \sigma^2_{\rho_{XY}} - 2\sigma_{\rho_{XY}} \rho_{XY} \quad [3]$$

Even if $D_{obs} = Y_i - X_i$ (i.e., X and Y are experimentally linked), the covariance with a similar but independently observed linked difference $D_{obs} = Y_i - X_i$ is the same. Hence Formula 3 is the appropriate numerator for Formula 2 in both cases.

The variance of D for the independent case is

$$\sigma^2_X + \sigma^2_Y - 2\sigma_{XY} \rho_{XY} \quad [4]$$

For the linked case, however, $D_{obs} = Y_i - X_i$, and the variance equals

$$\sigma^2_X + \sigma^2_Y - 2\sigma_{XY} \rho_{XY} \quad [5]$$

Hence for the linked case the reliability coefficient $\rho_{(Y_i - X_i)(Y_i - X_i)}$ is obtained by dividing Formula 3 by Formula 5:

$$\rho_{DD} = \frac{\sigma^2_{\rho_{XY}} + \sigma^2_{\rho_{XY}} - 2\sigma_{\rho_{XY}} \rho_{XY}}{\sigma^2_X + \sigma^2_Y - 2\sigma_{XY} \rho_{XY}} \quad [6]$$

whereas for the independent case the reliability coefficient $\rho_{(Y_i - X_i)(Y_i - X_i)}$ is Formula 3 divided by Formula 4:

$$\rho_{DD} = \frac{\sigma^2_{\rho_{XY}} + \sigma^2_{\rho_{XY}} - 2\sigma_{\rho_{XY}} \rho_{XY}}{\sigma^2_X + \sigma^2_Y - 2\sigma_{XY} \rho_{XY}} \quad [7]$$

Since $\rho_{XY} < \rho_{XY}$ in most instances, ρ_{DD} for the independent case will most likely be smaller than ρ_{DD} for the linked case. Both reliability coefficients are meaningful. They describe the correlation between differences observed according to different experimental designs. Distinctions like that between Equations 6 and 7 have to be made in considering the reliability of any composite, weighted or unweighted.

In the discussion that follows formulas are written in terms of the linked case; the substitution to fit the fully independent case will be obvious.

ESTIMATORS OF TRUE CHANGE

Primitive Formulas

Among the possible estimators of D_{obs} are three simple formulas.

Raw gain. The simplest formula is:

$$D = Y - X \quad [8] [1]$$

arrive at

$$\hat{X}_a = \frac{\rho_{XY} - \rho_{XY} \rho_{XY} X}{1 - \rho_{XY}^2} + \frac{\sigma_X(\rho_{XY} - \rho_{XY} \rho_{XY})}{\sigma_Y(1 - \rho_{XY}^2)} Y + \text{constant} \quad [16]$$

Then

$$\hat{D}_a = \hat{Y}_a - \hat{X}_a \quad [17]$$

where the estimates on the right-hand side come from Equation 16 and its analog for \hat{Y}_a . Expanding the expression we have

$$\begin{aligned} \hat{D}_a = & \frac{1}{1 - \rho_{XY}^2} \left[\frac{Y}{\sigma_Y} (\sigma_{YXY} - \rho_{XY} \sigma_{XY}) - \rho_{XY} \sigma_{XY} \right. \\ & + \sigma_X \rho_{XY} \sigma_{XX} - \sigma_X \rho_{XY} \left. - \frac{X}{\sigma_X} (\sigma_{XXY} - \rho_{XY} \sigma_{XY}) \right. \\ & \left. - \sigma_X \rho_{XY} \sigma_{XY} + \sigma_Y \rho_{XY} \sigma_{XY} - \sigma_Y \rho_{XY} \right] \\ & + \text{constant} \quad [17a] \end{aligned}$$

For independent observations, one substitutes ρ_{XY} wherever ρ_{XY} appears in Equation 17a; this is Lord's estimator of D_a . In an experiment, all calculations must be made within a treatment group.

Estimator 4 is as good or better than any of those listed ahead of it, giving a smaller mean square of $(\hat{D}_a - D_a)$ and a larger correlation between estimate \hat{D}_a and true score D_a . Ordinarily Estimator 3 is better than Estimator 2 and both are better than Estimator 1. Our main point in presenting Estimators 2 and 3 is to show the Lord estimator as an elaborated form of a more conventional estimator. This lays a base for the further refinement.

It may seem anachronistic in Formulas 11 and 16 to use a posttest score to "predict" a pretest score. But the logic is clear. Within a treatment, persons higher on the posttest than others having the same observed pretest score tend to be those for whom the true pretest score is higher than the observed score. The Y receives at least nominal weight in the regression equation when ρ_{XY} is not zero or one, and $\rho_{XY} < 1.00$. The weight given to Y increases with larger ρ_{XY} and smaller ρ_{XX} .

By taking group membership into account. Previous papers have implicitly assumed that all persons come from a single population, but often there are several distinct subgroups. These groups may be distinguished by demo-

graphic characteristics or by past experience, or they may be groups receiving different treatments between the X and Y observations. One could pool all groups and determine a single within-group value for each parameter in the equations above, but parameters calculated within subgroups will give a better estimate of X_a , Y_a , and D_a . There is a limit to how far subdivision of samples can profitably be carried, however.

Correlations and Regression Slopes for D_a

Sometimes an investigator wishes to know the correlation of D_a with another variable, say Q . He should note, then, that the correlation of D_a with Q is not a sound estimate of the correlation of D_a with Q . The correlation should be determined directly from the covariance of D_a with the second variable of interest. This covariance takes a form such as

$$\sigma_{D_a Q} = \sigma_{YQ} - \sigma_{XQ} = \sigma_{YQ} - \sigma_{XQ}$$

To get the correlation coefficient one divides by σ_Q and σ_{D_a} (not σ_{D_a}). Since D_a equals $Y_a - X_a$, its variance is given by Equation 3. All parameters must be those for the same group.

The investigator is often interested in the regression of D_a (or Y_a) on another variable. The slope of the regression of D_a on (e.g.) X_a is $\sigma_{D_a X_a} / \sigma_{X_a}^2$.

Attention must be paid to linkages. Let us write Q_i to indicate that the observation of Q is independent of X_i , Y_i , and Y_i .

$$\sigma_{D_a Q_i} = \sigma_{Y_i Q_i} - \sigma_{X_i Q_i} = \sigma_{Y_i Q_i} - \sigma_{X_i Q_i} \quad [18]$$

This cannot be determined from data where Q is linked to X or Y .

A variance-covariance algorithm. A simple computational routine can be suggested for problems of this character. One may form a variance-covariance matrix of observed values as in Figure 1. Linked and independent covariances are carefully distinguished. The matrix may be augmented as shown, adding rows and copying forward covariances. Reliability information is taken into account at certain points. Columns may be added to the matrix also, in a symmetric manner. Thus all entries in the X_a row can be copied into the X_a column. The full extension carried out in this way gives a square matrix, from which such values as $\sigma_{D_a Q_i}$ and $\sigma_{D_a}^2$ can be read out.

	X_a	Y_a	Q_a	Y_a	Q_a	X_a
Row 1 X_a	σ_{XX}	σ_{XY}	σ_{XQ}	σ_{XY}	σ_{XQ}	σ_{XX}
Row 2 Y_a	σ_{XY}	σ_{YY}	σ_{YQ}	σ_{YY}	σ_{YQ}	σ_{XY}
Row 3 Q_a	σ_{XQ}	σ_{YQ}	σ_{QQ}	σ_{YQ}	σ_{QQ}	σ_{XQ}
Row 4 Y_a	σ_{XY}	σ_{YY}	σ_{YQ}	σ_{YY}	σ_{YQ}	σ_{XY}
Row 5 Q_a	σ_{XQ}	σ_{YQ}	σ_{QQ}	σ_{YQ}	σ_{QQ}	σ_{XQ}
Row 6 X_a	σ_{XX}	σ_{XY}	σ_{XQ}	σ_{XY}	σ_{XQ}	σ_{XX}
Row 7 Y_a	σ_{XY}	σ_{YY}	σ_{YQ}	σ_{YY}	σ_{YQ}	σ_{XY}
Row 8 D_a						
Row 9 D_a						
Row 10 D_a						
Row 11 Q_a						

Copy independent covariances corresponding to Row 1

Copy independent covariances corresponding to Row 4

Fill in by subtracting Row 1 from Row 2

Fill in by subtracting Row 1 from Row 4

Fill in by subtracting Row 6 from Row 7

Copy independent covariances from Row 5

FIG. 1. Algorithm for constructing covariances and variances. (Covariances for linked observations are identified by the symbol σ_{XY} , and those for independent observations by σ_{XQ} . The broken line separates the original data from entries added later.)

A Better Estimate of True Change

Lord's formula uses only X and Y data, but we shall bring in two further categories of variables, W and Z . The W and X are Time-1 measures, but need not be simultaneous. Although we write W without vector notation, there are, in principle, any number of W variables that can be used singly or in combination. Our statements apply to any W or any weighted composite of the W .

A W might be any score describing the subject as he was prior to the treatment under study or W might be an index based on his life history. The scores Y and Z are posttreatment measures—again, not necessarily simultaneous. The Y might, for example, be a measure of performance at the end of training, and Z a retest a month later. Where we are examining a difference score rather than a change

score, no distinction between W and Z variables is needed.

The steps taken in going from Estimate 2 to Estimate 4 above can be extended to make use of W and Z information so as to reach an even better estimate of D_a .

The complete estimator. If W is univariate and there is no Z information,

$$\begin{aligned} \hat{X}_a = & \rho_{XY} X + \frac{\sigma_{XW} \sigma_{XZ}}{\sigma^2(W \cdot X)} (Y \cdot X) \\ & + \frac{\sigma_{XW} \sigma_{XZ}}{\sigma^2(W \cdot X \cdot Z)} (W \cdot X \cdot Z) + \text{constant} \quad [19] \end{aligned}$$

Here $W \cdot X \cdot Y$ is a partial variate, the deviation of W from the value predicted by the regression of W on X and Y .

If the W information is multivariate, a whole series of partial variates enters. The order of partialling is arbitrary; one might write terms

for X , W , X ; Y , W , X ; etc. Where there is Z information, one adds further terms, again employing partial variates; W , as well as X and Y , is partialled out. The estimation procedure must take into account any linkage between X , Y , W , and Z variables as in Equation 16.

A similar equation is written for \hat{Y}_a . Finally, one comes to an equation of the form

$$\hat{D}_a = \beta_1 X + \beta_2 Y + \beta_3 W + \beta_4 Z + \text{constant} \quad [20] [5]$$

Here $\beta_4 Z$ stands for a string of several terms of the form βW , if there are many W (likewise for $\beta_3 Z$). This estimator is superior to Formula 4, provided that the sample size is large enough to justify assigning a large number of weights. Where sample size is insufficient, the number of predictors must be held down, most likely by employing the first one or more principal components of the W set as predictors (likewise for Z).

When the problem becomes complicated, it is better to use efficient computing routines than to write out elaborate formulas. The within-treatment covariance matrix for X , Y , W , Z is written. Additional rows and columns for X_a and Y_a are formed as in Figure 1, with care to enter independent or linked covariances as required. The X_a column is subtracted from the Y_a column to form the D_a column. When the symmetric matrix is complete, one applies a multiple-regression program, treating entries in the D_a column as "test-criterion covariances" and the appropriate X , Y , W , and Z as predictors. If the observed scores have X and Y linked, for example, then X_a and Y_a are used as predictors, and covariances for Y_a are ignored.

Demographic information and information about experience can and should be considered as W variables. With a variable such as sex, one has a choice of entering it directly as a variable coded 1 and 0, or of performing a separate regression analysis for each sex. Both procedures regress the person's score toward the mean for his own sex rather than toward the mean for all cases. The second procedure allows for the possibility that the regression surface for males differs from that for females. Separate within-subgroup regressions would seemingly be preferred when samples are truly large.

The difficulty is that the argument can be repeated for every other noncontinuous variable and for all combinations of them. Indeed, it applies to continuous variables also; for example, a regression surface for more anxious children may differ from that for the less anxious. These remarks amount to entertaining the possibility of nonlinear relationships. While this possibility is real enough, one can rarely get usable estimates of nonlinear functions from samples of practical size (Birket, 1964; Goldberg, 1969). Hence, after one has divided the sample into a few salient subgroups, each having a suitably large size, the dummy-variable technique seems to be the only feasible way to handle a variety of nonquantitative information.

RESIDUALIZED GAINS

Developments parallel to those above lead to so-called "residual gains" or "base-free measures of change."

Alternative Estimators

The raw residual-gain score is defined by

$$D \cdot X = Y - E(Y) | X = Y - \hat{Y} - \beta_{Y \cdot X}(X - \bar{X}) \quad [21] [1]$$

We designate this Formula 1' to emphasize that it is comparable to Formula 1, the raw gain. If Y_a and X_a define the difference score, $\beta_{Y \cdot X}$ equals $\beta_{Y \cdot X} / \sigma_X$. If Y_a rather than Y_a is used, $\beta_{Y \cdot X}$ replaces $\beta_{Y \cdot X}$. The traditional definition given by Equation 21 is ambiguous; the residual ($Y_a - X_a$) $\cdot X_a$ is conceptually different from the residual ($Y_a - X_a$) $\cdot X_a$.

Residualizing removes from the posttest score, and hence from the gain, the portion that could have been predicted linearly from pretest status. One cannot argue that the residualized score is a "corrected" measure of gain, since in most studies the portion discarded includes some genuine and important change in the person. The residualized score is primarily a way of singling out individuals who changed more (or less) than expected.

True residual gain could be defined either as the expected value, over many observations on the same person, of $D \cdot X$ (defined in either of the two possible ways), or as the partial variate $D_a \cdot X_a$, the part of the true gain not predictable from true pretest status. $D_a \cdot X_a$ is more

likely to be of interest. Certainly if we intend to pick out superior learners or persons whose self-concept falls far below the self-ideal, we would like to base the discrimination on a true-score disparity. Likewise, if we have correlational questions—for example, Does anxiety predict overachievement?—the variable seems better specified by $D_a \cdot X_a$. Hence we consider only the definition:

$$D_a \cdot X_a = D_a - \beta_{D_a \cdot X_a} X_a + \text{constant} \\ = Y_a - \frac{\sigma_{Y_a X_a}}{\sigma_{X_a}} X_a + \text{constant} \quad [22]$$

No ambiguity arising from possible linkage of the observed X and Y enters this definition, though linkage must be considered in any estimation procedure.

Successive estimators of true residual change can be constructed on the same principles as before, but it will suffice to move directly to the estimator formed by the multiple-regression principle in the manner of Lord:

$$\hat{D}_a \cdot \hat{X}_a = \frac{\beta_{Y_a X_a} - \beta_{Y_a X_a} \beta_{X_a X_a}}{(\sigma_{X_a} - \beta_{Y_a X_a} \sigma_{X_a})} X_a + \text{constant} \quad [23] [4] \\ \times \left(Y_a - \frac{\sigma_{Y_a X_a}}{\sigma_{X_a}} X_a \right) \quad [23] [4]$$

This and other constants in this section are defined to make the mean of estimated residual gain equal zero. This estimate turns out to be proportional to the raw residual gain. (If the estimate is made from independent Y and X , $\beta_{Y_a X_a}$ replaces $\beta_{Y \cdot X}$.)

If there is W or Z information, a still better estimator is

$$\hat{D}_a \cdot \hat{X}_a = \beta' X + \beta' Y + \beta' W + \beta' Z + \text{constant} \quad [24] [5]$$

Only in rare cases will this be proportional to the observed $D \cdot X$.

The computational algorithm used before can be extended to obtain the desired weights for Formula 4' or 5'. Suppose we have filled out the square matrix for X , Y , X_a , Y_a , D_a . Then we can form the covariance of any variable with $D_a \cdot X_a$ very simply. For example,

$$\sigma_{D_a \cdot X_a} = \sigma_{D_a} \sigma_{X_a} - \frac{\sigma_{D_a X_a}}{\sigma_{X_a}} \sigma_{X_a} \sigma_{X_a} \quad [25]$$

Hence we may multiply every entry in the X_a

column of the matrix by $\sigma_{D_a X_a} / \sigma_{X_a}$, entering this in a column to one side. Subtracting the entry in this side column from the entry in the corresponding row of the D_a column gives a covariance to be entered in a $D_a \cdot X_a$ column. This column is ~~very~~ taken as a set of covariances of predictors with the criterion, and a multiple-regression program is applied.

The Tucker-Damarin-Messick Proposals

This analysis puts us in a position to review and clarify the rather puzzling paper entitled "A base-free measure of change" (Tucker, Damarin, & Messick, 1966; hereafter, referred to as TDM). They start much as we do by noting that the psychometrics of a change score applies to all kinds of difference scores. They suggest, as Lord did, that one should be most interested in the true difference score. They propose to divide this difference into two components, "one entirely dependent on the true score of the first or base-line test" and one "entirely independent of it." That is, they are interested in a true predicted gain and a true residual gain. As their abstract says, "equations for estimating both components are given." Since we shall recommend against use of their equations, we shall not go into details of their argument.

It might appear that TDM are concerned with estimating $E(D_a) | X_a$ and $D_a - E(D_a) | X_a (= D_a - X_a)$. The former, of course, is a linear function of X_a . TDM arrive at an equation (their Equation 26) that, in a form consistent with the present paper, is

$$\hat{D}_a \cdot \hat{X}_a = \left\{ \begin{array}{l} Y_a - \frac{\sigma_{Y_a X_a}}{\sigma_{X_a}} X_a \\ \text{or} \\ Y_a - \frac{\sigma_{Y_a X_a}}{\sigma_{X_a}} X_a \end{array} \right\} + \text{constant} \quad [26]$$

It is rather startling to find that this agrees with none of our formulas. It differs from Formula 1' in that X is replaced by X / σ_X . It differs by a further constant of proportionality from Formula 4'. The marked departure from Formula 4' is made the more puzzling by the favorable references of TDM to the Lord and McNemar papers and by their recommendation of Formula 4 for the gain score itself.

Personal communication with the authors verified that a failure of communication had occurred. As readers, we had given too little weight to one key phrase: The measures are "primarily intended for correlational work."

That is to say, TDM have no intention of interpreting "basefree" scores for individuals. Such scores are intended only as an intermediate step toward correlations. TDM offer an estimator that does not give the best least-squares estimate of individual "basefree" scores because they seek instead estimates that correlate zero with X_{10} .

Their intention is to determine correlations so as to learn what kinds of person show gains larger than would be predicted from the true pretest score. Correlating the estimated true residual gain from Equation 23 or 24 with another variable does not give the correlation TDM desire (unless $\rho_{\gamma\gamma}$ is 1.00). To understand this, consider for a moment the simple correlation of Y with Q . If we do not know Y scores but do know $\rho_{\gamma\gamma}$, we might estimate Y from X scores by the usual regression equation. Then $\rho_{\gamma Q}$ will not be a good approximation to $\rho_{\gamma Y}$; it will actually equal $\rho_{\gamma Q}$. In general, one who wants to interpret correlations, covariances, or regression slopes ought not to work from estimated scores. TDM intended to recommend that in such a line of research one calculate special-purpose scores by Formula 26 and then determine correlations. This appears to be unsound. TDM desire to obtain correlations with various Q 's and γ 's, which in their notation are the residual and predictable portions of true gain, respectively. But the TDM formulas generate fallible values g and w ; g equals γ plus an error. Obviously $\rho_{gQ} < \rho_{\gamma Q}$ and $\rho_{wQ} < \rho_{\gamma Q}$. As this is not explained by TDM, their paper is likely to mislead the reader. The confusion is reflected, and to some degree intensified, when Traub (1967) and Glass (1968) comment on the TDM formula.

While one might adapt the TDM statements to get $\rho_{\gamma Q}$, $\rho_{\gamma Y}$, etc., this is unnecessary. A straightforward manipulation of the matrix of observed covariances for X , Y , and Q (along with $\rho_{\gamma\gamma}$ and $\rho_{\gamma Y}$) yields the covariance of Q with $D_{\gamma} \cdot X_{10}$ (i.e., with γ). The $\sigma^2_{D_{\gamma} \cdot X_{10}} (= \sigma^2_{\gamma})$ needed to reach a correlation is simply the co-

* L. R. Tucker, F. Danarin, and S. Messick, personal communication, September 1968 and April 1969.

variance of D_{γ} with $D_{\gamma} \cdot X_{10}$ that we have already obtained. To get covariances for $D_{\gamma} - D_{\gamma} \cdot X_{10}$ (i.e., for $\hat{\gamma}$), one need only subtract column $D_{\gamma} \cdot X_{10}$ of the covariance matrix from column D_{γ} . And $\sigma^2_{D_{\gamma} - D_{\gamma} \cdot X_{10}} = \sigma^2_{\gamma} + \sigma^2_{\hat{\gamma}}$.

A MULTIVARIATE CONCERN

The older statement of the problem as "the measurement of gain" or of "residual gain" implies a special affinity between X and Y —they are seen as "the same variable" in some sense. But change is multivariate in nature.

Even when X and Y are determined by the same operation, they often do not represent the same psychological processes (Lord, 1958). At different stages of practice or development different processes contribute to performance of a task. Nor is this merely a matter of increased complexity; some processes drop out, some remain but contribute nothing to individual differences within an age group, some are replaced by qualitatively different processes. This does not rule out purely empirical studies of changes in the operationally defined variable. To assess such changes, even when one cannot describe them qualitatively, may be practically important. One must be careful not to fall into the trap of assuming that the changes are in a particular psychological attribute.

We may illustrate by referring to Fleishman's well-known studies of psychomotor scores at successive stages of practice (see Fleishman, 1966). On the first few trials, scores tend to correlate with cognitive measures; the usual speculation is that the high-cognitive subjects gain fastest because they most rapidly comprehend the instructions, display, strategy, etc. In a second stage, certain pretests of psychomotor ability correlate highest with scores, and the correlation of scores with cognitive pretests becomes rather small. One could easily conclude that cognitive ability is unimportant to learning in this second stage. But the drop in correlation (assuming no great rise in SD) demonstrates something more striking: that cognitive ability is negatively correlated with change from the first to second stage. And one can suggest a good reason. If bright persons

³ This matrix-extension procedure is entirely consistent with the TDM rationale. In fact, TDM tell us that their formula was arrived at by analytic treatment of this extended matrix.

catch on fastest, by some trial t they have completed what can be done intellectually with the task. On trials $t+1$ and later, their cognitive abilities produce no further gains. But by our hypothesis the dull have not comprehended fully by trial t , and therefore they have cognitive work left to do on trials $t+1, t+2, \dots$

During this stage, any gain in performance that comes from improved comprehension will be a gain by those low in cognitive aptitudes, hence those aptitudes have a negative correlation with gains. The Fleishman studies have not been processed in terms of gain scores, but it has been commonly said that they indicate cognitive abilities to be "important during the early stages of practice on psychomotor tasks," or the like. It seems more accurate to say that cognitive abilities make their contribution at different times for different persons and that gains at any one time are due to different processes for different persons.

Something similar is to be said about the relation of mental age (MA) at one age to subsequent gains in MA or achievement. A positive relation would result insofar as the high scorers understand new material better, or are more efficient learners. But there would also be a negative relation, insofar as the high scorers are those who have already mastered some highly valuable technique (e.g., mediation) that the low scorers have yet to master. As they restructure their behavior, the persons with low scores at the start of the period may make large gains—gains the high scorers had previously made. Positive and negative elements are probably both present, which should make us much less surprised than we have been by reports of near-zero correlation of MA (in a constant-age group) with subsequent gain in MA (Anderson, 1939; Bloom, 1964, p. 26 ff., 62 ff.).

We reduce emphasis on the special role of X as precursor of Y , and regard the whole $W \cdot X$ set as a vector describing the person's initial status. Then one may ask how Y varies as a function of the Time-1 data. To single out for intensive study persons who do better (or worse) than predicted, for example, it is wise to define expected outcome as the forecast of Y on the basis of all Time-1 information. Instead of $D_{\gamma} \cdot X_{10}$ ($= Y_{\gamma} \cdot X_{10}$) one would estimate $D_{\gamma} \cdot X_{10}$, W_{10} ($= Y_{10} \cdot X_{10}$, W_{10}). The machinery suggested above for partialling out X_{10} would

be used, but extended by partialling the W_{10} also out of D_{γ} . The entire set of W and X measures constitute the "base." There are optional targets for investigation: D_{γ} ; $D_{\gamma} \cdot X_{10}$; $D_{\gamma} \cdot X_{10}$; W_{10} ; etc.

Learning or growth is multidimensional; many measures could be taken at each point in time. To select one particular Y as somehow integrating a variety of subcriteria is to sacrifice information and possible insight. A person's change is better described by a vector of true scores W_{10} , X_{10} ; Y_{10} , Z_{10} . Each of these can be estimated by the methods used in Equations 19 and 20. One who wants to examine predicted and residual change will estimate \hat{Y}_{10} from W_{10} and X_{10} , and also from all variables together. He will obtain these scores:

$$\begin{aligned} \hat{Y}_{10} | WXYZ & \text{ (estimated true final status)} \\ \hat{Y}_{10} | W_{10} X_{10} & \text{ (predicted true final status)} \\ \text{Difference} & \text{ (estimate of unpredicted true} \\ & \text{residual)} \end{aligned}$$

The estimates of W_{10} and X_{10} come from W , X , Y , and Z . There would be estimates like those for Y for each Z or for orthogonal components of the Y_{10} , Z_{10} space.

We can rearrange the vectors W , X and Y , Z in a great variety of ways. Which target to choose can be decided only in the light of the purposes of the study.

PURPOSES OF ESTIMATING GAINS OR DIFFERENCES

Just why gains or differences are thought to be worth estimating can perhaps be inferred from the studies where estimates of some sort have been made in the past. The following aims may be noted:

1. To provide a dependent variable in an experiment on instruction, persuasion, or some other attempt to change behavior or beliefs.
2. To provide a measure of growth rate or learning rate that is to be predicted, as a way of answering the question, What kinds of persons grow (learn) fastest? Here, the change measure is a criterion variable in a correlational study.
3. To provide an indicator of deviant development, as a basis for identifying individuals to be given special treatment or to be studied clinically.
4. To provide an indicator of a construct that is thought to have significance in a certain

significant change, or to describe the magnitude of the effect. An estimate of true gain might appear to be pertinent. But it is not. For if one were to estimate D_w for each individual, and average, he would arrive back at the sample mean of observed gain. A significance test need only ask whether μ is reliably different from μ_0 . The difference in sample means for X and Y is the best available estimate of the mean D_w .

Criteria in Correlational Studies.

Correlational studies are often intended to investigate a question such as this: Among persons with a given pretest score, what attributes distinguish those who profit most from the treatment? This may seem to ask about ρ_{00} or ρ_{01} , or perhaps about $\rho_{(0-X)Y}$ or $\rho_{(0-X)W}$. It is more straightforward to ask about the regression of Y on X_w and W_w , the corresponding correlation for Y or Y_w , or related partial correlations. It appears that nothing is gained by referring to change measures in this context. The relationships of true scores can be investigated without estimating true scores for individuals.

Selecting Individuals on the Basis of Gain or Difference Scores.

Many who calculate difference scores are interested in making decisions about individuals—identifying underachievers for clinical attention or fast learners for special opportunities, for example. One can scarcely defend selecting such individuals on a raw-gain or raw-difference score, especially as these scores tend to show a spurious advantage for persons low on X . Selecting cases whose estimated Y_w is higher than that of others with similar X_w and W_w seems more sensible. To do this, regression equations should be called into play. That is, one selects persons for whom Y_w is much larger (or smaller) than Y_w than $W_w X_w$.

The persons with positive deviations are those who did better than predicted. This means either that they started with some valuable attribute the W and X variables did not encompass, that their pretest-true scores are underestimated or their posttest scores are overestimated, or that their success on Y was an accidental effect arising from some tactically adopted during learning or some sequence of lucky trials. It is very hard to dispose

involve calculating residual gain scores for individuals.

Comparison of treatment groups not formed at random. When treatments are applied to groups differentiated by a nonrandom process, the X_w distributions within the subpopulations represented by the groups are generally not the same. Consequently, the same observed X score implies a different level of true pretest ability, depending on the group.

If analysis of covariance is to be made, it is advisable to regress the covariate toward the mean of the treatment group before entering it in the analysis (Lord, 1960). If there is W information as well as X , it also contributes to the estimate. So does Y and Z information. Here is a paradox: A proposal to use the posttest score to estimate the pretest true score which will then be used to adjust posttest scores! The crucial point is that the estimator of the covariate is determined from within-group data. Since the estimate of any linear group mean as that function of X and W , the procedure does not introduce bias nor does it reduce any effect truly attributable to the treatment.

Application of analysis of covariance to studies where initial assignment was nonrandom, which was widely recommended 10 years ago, is now in bad repute. Even the elaborate technique just suggested is no more than a palliative. If the treatment groups differed systematically at the start of the experiment with respect to any relevant characteristic other than the covariate, even a perfect measure of the covariate cannot remove the confounding. To quote Lord (1967), "there simply is no logical or statistical procedure that can be counted on to make proper allowances for uncontrolled preexisting differences between groups [p. 305]." And Meehl (1970) calls such corrections "inherently fallacious [in press]."

The findings of the study can be usefully summarized by calculating within-group regression functions relating Y_w to X_w , W_w , using the covariance matrix for true scores. What cannot be done is to "compare treatment effects."

One-group designs. A third kind of experiment is the simple one-group study where one wishes to learn whether a treatment produces

theoretical network. The indicator may be used as an independent variable, covariate, dependent variable, etc. An example is the interference score needed on the Stroop Color Word Test to represent the decline in reading rate when color names are presented in colors incongruent with the names.

Much of the confusion in the literature arises from a failure to distinguish these purposes and to match distinct methodological recommendations to them.

Gains as a Consequence of Treatments

There appears to be no need to use measures of change as dependent variables and no virtue in using them. If one is testing the null hypothesis that two treatments have the same effect, the essential question is whether posttest Y_w scores vary from group to group. Assuming that errors of measurement of Y are random, Y is an entirely suitable dependent variable.

The randomized experiment. Suppose that cases are assigned to treatments in a random or stratified-random manner. The X scores will vary within groups. An analysis of covariance to take this variation into account is advantageous so long as ρ_{XY} is large. (If $\rho < 0.4$, blockage is probably to be preferred, according to Elashoff, 1969.) The usual adjustment estimates the Y scores expected under the null hypothesis and then expresses each observed Y as a deviation from the estimate. Ordinarily it is desirable to base the adjustment, not on X , but on whatever linear combination of X and W best predicts Y within groups.

Where within-treatment regressions are linear but significantly different in slope, the difference between effects of treatments depends on the level of X . The "main effect" is not interpretable. The most meaningful report consists of regression functions for Y on the X_w , W_w space, computed with the aid of the covariance matrix for true scores within each group in turn.

Nowhere in this section have we made use of a change score. We consider it likely that a change will vary systematically with X_w . Where this is the case, the essential result is a regression function, not a mean gain. The adjusted Y score of the significance test is a sort of residual gain, but the procedure does not

of the hypothesis that these unexpected gains were fortuitous.

Here, the focus of attention is on an estimated residual gain: not $D \cdot X$, not $D_w \cdot X_w$, but $\hat{Y}_w - X_w \hat{W}_w$ or, what is equivalent, $D_w \cdot X_w \hat{W}_w$. Where X alone is available as a predictor, the raw residual gain selects the same persons as Formula 23 does. But Formula 24 is to be preferred.

It is possible of course, given before-and-after scores on the same instrument, to estimate true gains of individuals and to identify those who did and did not gain. But to what purpose? This has no clear bearing on decisions about the future of these persons, and the decision rule for fresh cases is to be inferred from the regression surface.

Differences and Gain Scores as Constructs.

One of the most common uses of difference scores is to operationalize a concept: For example, self-satisfaction is sometimes defined as the difference between the rating of self and ideal-self on an esteem scale. One might likewise think of a gain score as reflecting "learning ability" on a certain task. Operational definitions will often take the form of linear combinations of operations.

But there is little a priori basis for pinning one's faith on $Y_w - X_w$ as distinct from the more general $Y_w - eX_w$. Just what weight to assign the "correcting" variable is an empirical question. To arbitrarily confine interest to D_w (which means that a is fixed at 1.00) is to rule out possible discoveries. This argues, then, for discovering what function of Y_w and X_w has the strongest relationships with variables that should connect with the construct.

The claim that an index has validity as a measure of some construct carries a considerable burden of proof. There is little reason to believe and much empirical reason to disbelieve the contention that some arbitrarily weighted function of two variables will properly define a construct. More often, the profitable strategy is to use the two variables separately in the analysis so as to allow for complex relationships.

One example of an "obvious" but questionable use of a subtractive correction is provided by a study in which skin conductance is a variable. At the start of the experiment a "base-

line" measure of the subject's galvanic skin response is taken. Then stress is applied and a second measure is taken. During a rest period the subject receives a drug or a placebo. Stress is again applied and a third measure taken. Call the measures, in order, W , X , and Y . Simple correction would use $Y - W$ as dependent variable and $X - W$ as covariate. We, however, would prefer to use X and W as separate covariates, with Y as dependent variable. This should give a more precise analysis when W is unreliable. (As suggested earlier, it would generally be still better to use X , W , XY and WXY as covariates.)

SUMMARY

Where true scores for individuals are desired, multiple regression procedures outlined herein make use of more information than do procedures hitherto advanced. There seems to be no occasion to estimate true gain scores. In the experiment where treatment groups are formed nonrandomly, estimates of true scores on the covariate can reduce the resulting bias.

Where individuals who have exceptionally high or low residual gains are to be identified, the raw residual gain serves as well as the alternate formulas hitherto advanced. To estimate the individual's true residual gain, however, a superior formula is available.

Where correlations and regression functions relating true gains or true residual gains to other variables are desired, a calculating routine is available that makes it unnecessary to estimate gain scores for individuals.

It appears that investigators who ask questions regarding gain scores would ordinarily be better advised to frame their questions in other ways.

REFERENCES

- ANDERSON, J. E. The limitations of infant and preschool tests in the measurement of intelligence. *Journal of Psychology*, 1939, 8, 351-379.
- BLOOM, B. S. *Stability and change in human characteristics*. New York: Wiley, 1964.
- BURKETT, G. R. A study of reduced rank models for multiple prediction. *Psychometric Monographs*, 1964, No. 12.
- DUBOIS, P. H. *Multivariate correlational analysis*. New

- York: Harper, 1957.
- ELASSOFF, J. D. Analysis of covariance: a delicate instrument. *American Educational Research Journal*, 1969, 6, 383-402.
- FLINZMAN, E. A. Human abilities and the acquisition of skill. In E. A. Bilodeau (Ed.), *Acquisition of skill*. New York: Academic Press, 1966.
- GLASS, G. V. Response to Traub's "Note on the reliability of residual change scores." *Journal of Educational Measurement*, 1968, 5, 265-267.
- GLASS, G. V., CRONBACH, L. J., & RAJARATNAM, N. Generalizability of scores influenced by multiple sources of variance. *Psychometrika*, 1965, 30, 395-418.
- GORIUNOV, L. R. The search for configurational relationships in personality assessment: the diagnosis of psychosis vs. neurosis from the MMPI. *Multivariate Behavioral Research*, 1969, 4, 523-536.
- HAKKIOVIST, K. Relative changes in intelligence from 13 to 18. *Scandinavian Journal of Psychology*, 1968, 9, 59-82.
- HARRIS, C. W. (Ed.). *Problems in measuring change*. Madison: University of Wisconsin Press, 1963.
- LOBO, F. M. The measurement of growth. *Educational and Psychological Measurement*, 1936, 16, 421-437.
- See also Errata, *ibid.*, 1957, 17, 452.
- LOBO, F. M. Further problems in the measurement of growth. *Educational and Psychological Measurement*, 1958, 18, 437-454.
- LOBO, F. M. Large-sample covariance analysis when the control variable is fallible. *Journal of the American Statistical Association*, 1960, 55, 309-321.
- LOBO, F. M. Elementary models for measuring change. In C. W. Harris (Ed.), *Problems in measuring change*. Madison: University of Wisconsin Press, 1963.
- LOBO, F. M. A paradox in the interpretation of group comparisons. *Psychological Bulletin*, 1967, 68, 304-305.
- MCNEAM, O. On growth measurement. *Educational and Psychological Measurement*, 1958, 18, 47-55.
- MEZENT, P. E. Nuisance variables and the ex post facto design. In M. Radner & S. Winokur (Eds.), *Minnesota studies in the philosophy of science*, IV. Minneapolis: University of Minnesota Press, 1970, in press.
- STRASZY, J. C. General and special formulas for reliability of differences. *Journal of Educational Measurement*, 1967, 4, 249-252.
- TRACY, R. E. A note on the reliability of residual change scores. *Journal of Educational Measurement*, 1967, 4, 253-256.
- TRUMBULL, H. C., & CRONBACH, L. J. A practical procedure for the rigorous interpretation of test-retest scores in terms of pupil growth. *Journal of Educational Research*, 1943, 35, 481-488.
- TRUCKER, L. R., DAMARIS, F., & MESSICK, S. A baseline measure of change. *Psychometrika*, 1966, 31, 457-473.

(Received June 24, 1969)