

ESTIMATION OF LATENT ABILITY AND ITEM PARAMETERS WHEN THERE ARE OMITTED RESPONSES*

FREDERIC M. LORD

EDUCATIONAL TESTING SERVICE

Omitted items cannot properly be treated as wrong when estimating ability and item parameters. A convenient method for utilizing the information provided by omissions is presented. Theoretical and empirical justifications are presented for the estimates obtained by the new method.

At the time the likelihood equations for item characteristic curve (icc) theory were written down [Lord, 1953], there seemed to be three major obstacles to practical applications:

1. Solution of the likelihood equations for data of real interest did not seem practicable from a computational point of view [Torgerson, 1958, p. 388].
2. Icc theory dealt with unspeeded tests, whereas almost all standard tests are administered with a time limit that prevents some examinees from finishing.
3. Icc theory was first developed for dichotomous items. Typical test answer sheets, however, carry at least three distinct types of examinee response: correct response, wrong response, no response ("omits"). These three types of response often receive different scoring weights (for example, 1, $-\frac{1}{4}$, and 0, respectively), but even if omitted responses are scored as wrong, they can not reasonably be treated as wrong in the likelihood equations.

Solution of equations. Numerical solutions to the likelihood equations have now successfully been obtained for large data sets [Lord, 1968; Bock & Lieberman, 1970; Bock, 1972]. Even the maximum likelihood estimate of the parameter representing the icc lower asymptote, sometimes incorrectly called the chance level, can now be obtained by maximum likelihood [Wingersky & Lord, 1973] whenever there is enough data to determine this part of the curve.

Speededness. A time-limit test is, for some examinees, partly a measure of something called speed, which is quite distinct from the ability measured

* Research reported in this paper has been supported by grant GB-32781X from National Science Foundation.

by the *power score* that would be obtained if the test were administered without time limit. Icc theory is not presently equipped to deal with the speed dimension explicitly, but the theory can still be used to analyze answer sheets obtained in timed test administrations. To do this requires the assumption that examinees answer test questions in order. For each examinee, the items following his last recorded response (this item is called the *last item attempted*) are ignored throughout the estimation process. (In practice, examinees answering less than a third, say, of the n items may be omitted from all analyses.) Thus examinee ability θ is estimated from his responses to items presumed reached and item parameters are estimated from the responses of examinees who presumably reached the item.

This does not complicate the likelihood equations or the process of solving them. A key property of icc theory is that item parameters do not depend on the group of examinees tested, within reasonable limits; and that examinee ability (θ) does not depend on the items administered, assuming that all items measure the same psychological dimension. Thus, ignoring various examinees and various items ought not to have serious effects.

It is true that a θ estimated in this way will approximate the examinee's ability under power conditions only if his responses to the items actually attempted would have been the same in the absence of a time limit. Regardless of this, a θ estimated in this way appropriately reflects the examinee's effective ability level under the timed conditions actually provided.

If an examinee does not reach an item because of the time limit, this fact contains no usable information for inferring his ability level θ under the icc model. The one-dimensional icc model considered here provides no way to make use of any relationship (ordinarily curvilinear) that may exist between speed of response and θ . The term *omitted response* or simply *omit* will hereafter refer only to items presumed reached, not to items after the "last item attempted."

Omitted items. The present paper proposes and discusses a method for the effective use of the information represented by omitted responses when estimating ability and item parameters. Another method has been proposed and used by Bock [1972]. With the implementation of adequate methods for dealing with omitted responses, the application of icc theory to typical testing data is now effective and practical.

Omitted Responses

The meaning of an omitted response varies depending on the type of item and how the test is scored. It will be assumed throughout this paper that the items are multiple choice. If the score is the number of right answers, to be denoted by x , then the examinee who omits any item is acting against his own best interests. We will not consider such cases here.

The only common alternative to the number-right score is the *formula*

score containing a penalty for wrong answers. If A is the number of alternative responses provided for the test item, the usual formula score is

$$(1) \quad y_a = x_a - \frac{w_a}{A - 1},$$

where w_a is the number of wrong responses given by examinee a .

We will assume hereafter that (1) is used; also that examinees wish to maximize their expected scores and that they are fully informed about their best strategy for doing this. Under these conditions, an examinee should omit an item only if he believes his chance of success on the item is no greater than $c \equiv 1/A$. On the other hand, since the item has A alternative responses, his chance of success should not be less than c , since he can always do this well by strictly random guessing. Following this reasoning, we will assume hereafter that if an examinee were required to respond to a long series of A -choice items that he had omitted, his proportion of correct answers would be c . [Slakter (1968) presents empirical evidence that examinees omit more items than they should according to this assumption; however, his examinees were not explicitly instructed as to their best strategy.]

The ICC Model

In item characteristic curve (icc) theory, the probability that examinee a will answer item i correctly may be denoted by $P_{ia} \equiv P_i(\theta_a)$, here assumed to be an increasing, twice differentiable function of his ability θ_a . It might seem from the preceding paragraph that $P_i(\theta_a) = c$ whenever examinee a is required to answer an item i that he previously had omitted. This cannot be correct, however. If examinee b omits the same item, we would have $P_i(\theta_b) = c$, from which it would follow that $\theta_a = \theta_b$. Since two examinees who omit the same item may be at very different ability levels, it is clear that the probability c is a different kind of probability than P_{ia} .

It might seem natural to think of P_{ia} as the relative frequency of correct answers when item i is repeatedly administered to examinee a under some hypothetical conditions requiring him to forget his previous responses. This interpretation of P_{ia} is considered in detail by Meredith [1965]. We cannot use this interpretation here (nor in most practical work). Examinee a might know the answer to item $i = 1$ and so have a probability of 1 of answering it correctly. At the same time, he might be misinformed about item $i = 2$ and so have a probability of 0 of answering it correctly. At the same time, examinee b might have probability 0 of answering item 1 correctly and probability 1 of answering item 2 correctly. If the items measure the same trait, the four equations $P_1(\theta_a) = P_2(\theta_b) = 1$ and $P_2(\theta_a) = P_1(\theta_b) = 0$ are very difficult to reconcile.

P_{ia} is most simply interpreted as the probability that examinee a will give the right answer to a randomly chosen item having the same icc as

item i . An alternative interpretation is that $P_i(\theta_a)$ is the probability that item i will be answered correctly by a randomly chosen examinee of ability level $\theta = \theta_a$. These two interpretations will usually be assumed to hold simultaneously.

These interpretations tell us nothing about the probability that a specified examinee will answer a specified item correctly. It is this last probability that would equal c if an examinee were required to answer an item he has omitted.

The Likelihood Function

Let the response of examinee a to item i be denoted by u_{ia} . For a correct response, let $u_{ia} = 1$; for an incorrect response, let $u_{ia} = 0$. When examinee a answers a test composed of n items, the u_{ia} are assumed independent (assumption of *local independence*). If he does not omit any items, the likelihood function for his responses can be written

$$(2) \quad L_a(u_{1a}, u_{2a}, \dots, u_{na} \mid \theta_a) = \prod_{i=1}^n P_{ia}^{u_{ia}} Q_{ia}^{1-u_{ia}}$$

where $Q_{ia} \equiv 1 - P_{ia}$. If the *item parameters* in P_{ia} that characterize item i are known (approximately, from pretesting), then the maximum likelihood estimate $\hat{\theta}_a$ of the examinee's ability can be obtained from this likelihood function by standard procedures.

If the item parameters are unknown, they can be estimated at the same time as θ_a from the responses of many (preferably two or three thousand) examinees. In this case, the likelihood function is

$$(3) \quad L(U \mid \theta) = \prod_{a=1}^N \prod_{i=1}^n P_{ia}^{u_{ia}} Q_{ia}^{1-u_{ia}},$$

where U is the matrix $\|u_{ia}\|$ and θ is the vector $\{\theta_1, \theta_2, \dots, \theta_N\}$. The maximum likelihood estimates of θ and of the item parameters can actually be obtained in practice from this likelihood function by standard procedures (somewhat surprisingly, in view of the very large number of parameters to be estimated). The practical effectiveness of this procedure is being demonstrated in work with real data [for example, Lord, 1970] despite the (presumably temporary) lack of a rigorous proof that the maximum likelihood estimates are consistent.

If the examinee omits certain items, it might seem that one could simply omit these items altogether from (2) or (3) and proceed as before. This cannot be right, however, since the fact that the examinee omitted certain items carries the important information that he did not know the answers to these items—that his chance of success was roughly only c on each. We cannot afford to ignore this information. If we did, an examinee could obtain

as high a $\hat{\theta}$ as he wished, simply by omitting questions whenever he was not completely sure of the correct answer.

One way to deal with this situation would be for the psychometrician to replace each omitted response by a response drawn at random with probability of success c . After all, this is just what some examinees do in actual practice, instead of omitting items. According to the model, the likelihood function (2) or (3) will hold for the "filled-in" data so obtained.

Although this procedure should yield consistent estimates (as $n \rightarrow \infty$), it is objectionable from two related points of view. From the examinee's point of view, it is unfair to saddle him with a possibly unfortunate set of random responses. From the statistician's point of view, the procedure degrades the data by introducing random error; it can only increase error variance, it cannot possibly be truly beneficial.

It would be desirable to replace (2) or (3) by a likelihood function that includes provision for omitted responses. Such a function, however, would depend in part on the true probability that examinee a will omit a randomly chosen item having the same item parameters as item i ($i = 1, 2, \dots, n$). This true probability would be a function, similar to $P_i(\theta_a)$ but not the same, depending not only on θ_a and on certain characteristics of the item, but also on a new trait of the examinee representing his willingness to omit items. Even after simplifying assumptions, there would be at least one new examinee parameter and one new item parameter to estimate, considerably complicating the already complicated and expensive estimation procedure.

New Estimation Procedure

The following estimation procedure has been used on several large sets of data, apparently with great success, as briefly indicated in the next section. The likelihood function (2) is replaced by

$$(4) \quad L_a^*(v_{1a}, v_{2a}, \dots, v_{na} \mid \theta_a) = \prod_{i=1}^n P_{ia}^{v_{ia}} Q_{ia}^{1-v_{ia}}$$

where $v_{ia} = 1, 0$, or c according to whether the response is right, wrong, or omitted, respectively. These values of v_{ia} are the item scores used in a familiar method of formula scoring. The score $\Sigma_i v_{ia}$ is perfectly correlated with the formula score given by (1).

Since for the present the item parameters must be estimated at the same time as θ_a , the responses of many (preferably two or three thousand) examinees are analyzed simultaneously, so that (4) is replaced by

$$(5) \quad L^*(V \mid \theta) = \prod_{a=1}^N \prod_{i=1}^n P_{ia}^{v_{ia}} Q_{ia}^{1-v_{ia}},$$

where V is the matrix $||v_{ia}||$. The new procedure, as actually used, finds the value of θ_a for each of the N examinees and the values of three item parameters

for each of the n items so as to maximize (5). These values are taken as the parameter estimates desired.

Since (4) and (5) are not likelihood functions, these estimates are not maximum likelihood estimates. The estimates of θ from (4) or (5) will be denoted by θ^* . It is shown in the Appendix that in the case of (4), the θ^* converge for large n to the same values as do the $\hat{\theta}$ obtained from (2) after omits have been replaced by random responses. Moreover, if there are omits, the sampling error of θ^* for large n is less than the sampling error of the maximum likelihood estimate obtained from the filled-in data.

Discussion of Assumptions

Many people prefer number-right scoring to the formula scoring considered here. Some common misinterpretations will be avoided by pointing out that certain assumptions are *not* made here.

1. There is no assumption here that formula scores are superior to number-right scores. If number-right scores were used, the problem considered here should not arise, since in that case examinees should not omit any items at all.
2. There is no assumption here that examinees guess at random when they do not know the answer to an item. On the contrary, many of the item characteristic curves found so far in the analysis of nationally used tests show that very low-ability examinees tend to do less well on difficult items than they would have done if they had responded at random. This situation presumably arises because certain of the possible item responses have been cleverly made so attractive that low-level examinees tend to choose them in preference to the correct answer.
3. The model given here is consistent with the obvious fact that examinees use misinformation and partial information in answering items. For a few items in nationally administered tests, icc have been found that never go below $P_i(\theta) = .30$ or $P_i(\theta) = .40$, regardless of θ level. This suggests that on some items even very low-level examinees may be able to rule out two or three of the possible item responses as incorrect.

It is assumed that the probability of a correct answer would be $c \equiv 1/A$ if an examinee were required to respond to the A -choice items he has omitted. As explained previously, this assumption is made because an examinee who wants to maximize his expected score should never omit an item if he can do better than choose among the A responses at random. Adequate directions to the examinee should convince him of this.

Empirical Results

Several questions can be raised about the various parameter estimates that have been discussed. In the first place, many people distrust the as-

sumptions of the ice model, particularly the assumption that there is only one dimension θ underlying the test. The best way to resolve this question is not to try to prove that the assumptions hold for a particular set of data (they will never hold exactly), but to show that the parameter estimates obtained provide a useful and effective summary of the data, capable of predicting new sets of data not yet observed. The main purpose of this section is to show just that, insofar as possible with the limited investigations made to date.

Maximum likelihood estimates obtained from (3) are open to a further objection—there exists no rigorous proof that these estimates are consistent. A related but distinct problem is that it may seem hard to believe that several thousand parameters can really be successfully estimated simultaneously. It would be valuable to have a mathematical proof of the asymptotic properties of the estimates. Any final answer to both questions, however, must come by demonstrating the usefulness of estimates obtainable in practice from samples of reasonable size.

Estimates obtained from (4) and (5) are open to further objections. Equations (4) and (5) are not likelihood functions. No clear statistical justification has been given for choosing these functions to be maximized (there seems to be an interesting and unanswered question in statistical inference here). Justification of the estimates obtained is given in the Appendix, but it still is important to show the practical effectiveness of actual estimates.

Data and Method

Estimates obtained by the new method (5) have been obtained for three sets of real data:

V. 2926 examinees tested on a 90-item verbal aptitude test.

v. 994 examinees tested on a parallel form of the same 90-item verbal aptitude test.

M. 2946 examinees tested on an 85-item arithmetic reasoning test.

m. same data as *M* leaving out the first 25 items (*m* is an actual test, *M* is a composite of two tests).

These tests have an unusually wide range of item difficulty. For each of tests *V*, *v*, and *M*, there was an average of about 9 omitted items per examinee. All items were five-choice items.

In order to permit comparisons, the answer sheets for tests *V* and *M* were reanalyzed after replacing omitted responses ("omits") by random responses, obtaining estimates by maximizing the likelihood (3). The estimates so obtained will be referred to as the MLE. Note that "not reached" items are not treated as "omits" in this study.

All estimation procedures were based on the three-parameter logistic model [Birnbbaum, 1968]. Investigations 1 and 2 below illustrate the effectiveness of (5) for certain applications. Investigations 3 and 4 illustrate the

effectiveness of the MLE from (3). Investigation 5 evaluates the degree of agreement between estimates from (3) and estimates from (5). Investigation 6 suggests that the θ^* actually are better estimates of θ than are the $\hat{\theta}$. Investigation 7 generates artificial data and reports the relation of parameter estimates from (5) to the true parameters. Additional, obviously desirable investigations have not been carried out because of the considerable computer costs.

1. *Relation of estimated ability to test score.* Figure 1 shows for test m the relation of the ability estimates θ_a^* obtained from (5) to formula score on the total test. Under the probabilistic icc model, scatter of scores about the test characteristic curve (the regression of test score on ability) is to be expected, because of sampling fluctuations provided for in the model. The correlation ratio of formula score on θ^* was .978.

For test v , a correlation ratio of .982 was computed, but coarse grouping

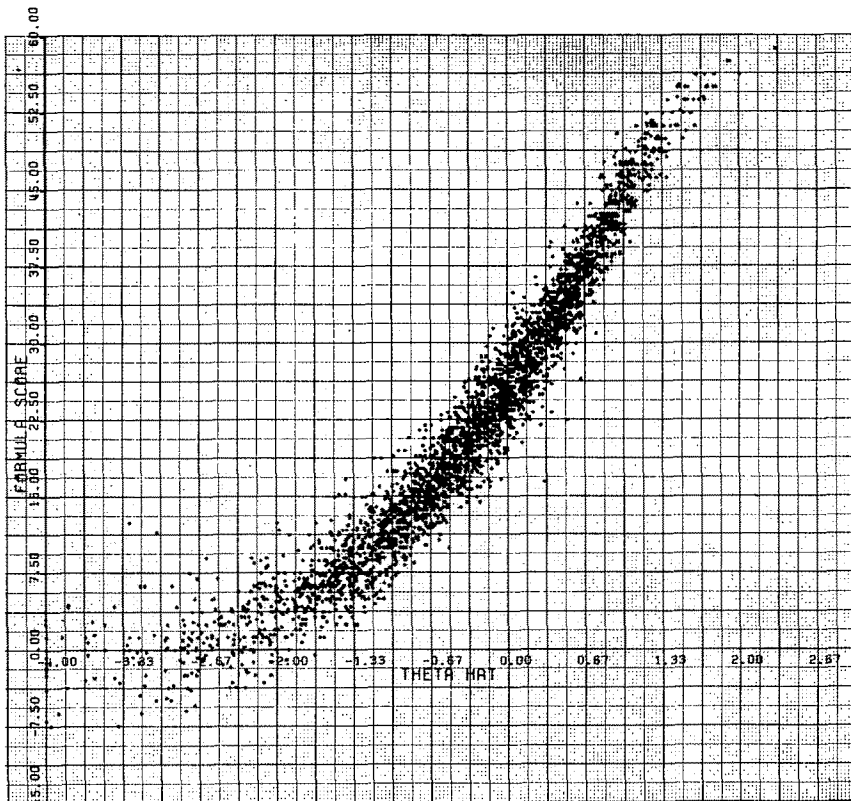


FIGURE 1

Relation between estimated ability (θ^*) and formula score.

of formula scores makes this somewhat too low. Tests V and M have not been analyzed in this way.

2. *Predicting test score from estimated ability.* For given θ_a , the expected number-right score of examinee a is

$$(6) \quad \varepsilon x_a = \sum^{n_a} P_i(\theta_a)$$

where \sum^{n_a} represents summation over all items actually answered. The correlation between $\sum^{n_a} P_i(\theta_a^*)$ and number-right score corrected for omits was obtained, separately for tests V and m . Here $P_i(\theta_a^*)$ represents $P_i(\theta_a)$ with parameter estimates from (5) substituted for the unknown values. The number-right score corrected for omits is $x_a + o_a/A$, where o_a is the number of omitted items. For 60-item m , the correlation was .981; for 90-item V , .992. Tests v and M have not been analyzed in this way.

These high correlations show that the estimated parameters summarize the data on the examinees' answer sheets very effectively. The value of this for practical applications is illustrated by Lord [1973].

Again, some scatter of scores about their expected value, due to sampling fluctuations, is provided for in the model. The scatter is less than it would be if true parameters had been used instead of estimated parameters. The reason is that chance irregularities in the data are to some extent fitted in the course of the estimation process. Cross validation procedures, using a second random sample of data, could be used to eliminate this.

3. *Comparing estimates of the distribution of ability.* The histogram in Figure 2 shows for the 1807 examinees who answered the last item in test m the actual frequency distribution of the $\hat{\theta}$ obtained from the filled-in data by (3). The smooth curve shows an estimate $\hat{h}(\theta)$ of the frequency distribution of θ (not $\hat{\theta}$) obtained by the method outlined below. The two distributions are obtained in very different ways, under totally different sets of assumptions, as detailed in Lord [1970], where θ were estimated by the two different methods and compared.

In order to obtain $\hat{h}(\theta)$, the frequency distribution $g(\xi)$ of true score ξ was estimated from the observed distribution of number-right scores under the compound-binomial error model [Lord, 1969; Wingersky, Lees, Lennon, & Lord, 1969]. By definition, true score is expected observed score, so by (6)

$$(7) \quad \xi_a = \sum^n P_i(\theta_a).$$

When the item parameters are known, this equation defines ξ as a function of θ , or θ as a function of ξ . Thus the distribution of ability $h(\theta)$ was estimated for the group of 1807 examinees by first estimating the distribution $g(\xi)$ from the distribution of their number-right scores on the filled-in answer sheets and then transforming this distribution to that of θ by the functional

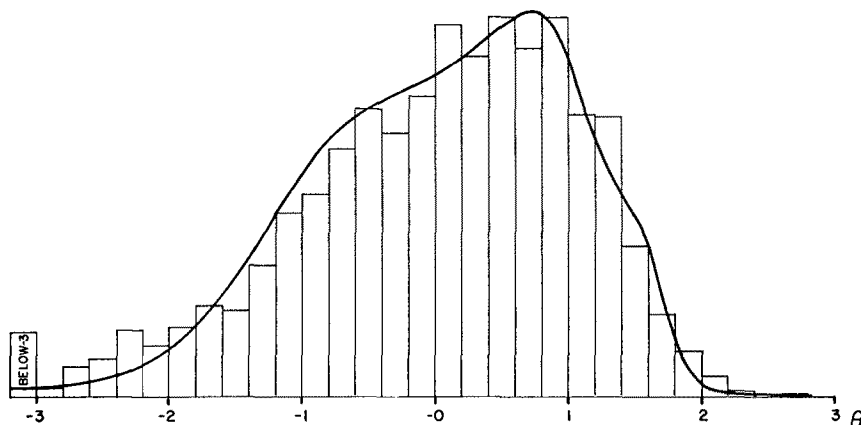


FIGURE 2

Distribution of estimated θ (histogram) and estimated distribution of θ (curve).

relationship (7) with estimated item parameters substituted for their unknown true values.

The two distributions in Figure 2 agree remarkably well. Even the discrepancies occur where they should. The $\hat{\theta}$ properly show a slightly more dispersed distribution than the θ (represented by the smooth curve), since the $\hat{\theta}$ contain errors of estimation. Because m is a difficult test, these errors are quite large for low ability examinees, as discussed in an earlier section.

Very similar results were obtained for test V . Tests M and v were not analyzed in this way. In view of the very different assumptions made, the correspondence shown in Figure 2, also in the unpublished figure for test V , is strong evidence for the meaningfulness and practical usefulness of the models and estimation procedures used to obtain these results.

4. *Estimated frequency distribution of number-right scores.* According to the icc model, the probability generating function [Kendall & Stuart, 1958, section 1.37] for the frequency distribution of number-right scores is

$$(8) \quad \sum_{x=0}^n \phi(x) t^x \equiv \int_{-\infty}^{\infty} \prod_{i=1}^n [Q_i(\theta) + tP_i(\theta)] h(\theta) d\theta.$$

Using the estimated $h(\theta)$ in Figure 2 and using estimated item parameters in $P_i(\theta)$ and $Q_i(\theta)$, the frequency distribution of number-right scores was estimated from (8). The resulting $\hat{\phi}(x)$ agreed to at least three decimal places with the $\hat{\phi}(x)$ obtained under the compound-binomial error model, which, it should be noted, makes no use at all of any item parameters. The latter $\hat{\phi}(x)$ agree well with the actually observed distribution of scores: when the frequencies are grouped into 38 class intervals and 4 parameters are estimated from the data, the calculated chi square under the compound-binomial

model is near the 30-th percentile of the χ^2 distribution for 34 degrees of freedom.

Results for test V were the same, except that the chi square calculated under the compound binomial model was at the 67-th percentile with 33 degrees of freedom.

5. Correlation of new estimates and MLE. For test M , correlations were calculated between item parameter estimates obtained by the two methods (eqs. 3 and 5). For the lower asymptote (c_i) of the icc, the correlation was .990. For the discriminating power (a_i) the correlation was .995. For the difficulty parameter (b_i) the correlation was .9996. These results show that for purposes of estimating item parameters, the new estimation method yields results virtually equivalent to a maximum likelihood procedure based on the usual icc model and filled-in observations.

Tests v , V , m , and M are all difficult for low ability examinees. As a result, the ability parameter θ cannot be estimated accurately at the lowest levels—one cannot effectively distinguish between $\theta_a = -5$ and $\theta_a = -500$. When the 56 examinees* with estimated θ 's below 3.0 are omitted, the correlation between estimated θ 's obtained from (3) and (5) for test M is .997.

There were on the average about 6 omits per examinee for test m . In view of the high values of all four correlations found, similar correlations were not computed for other tests.

6. Likelihoods. For tests M and V , a new set of data ("CV data") was set up for cross-validation purposes by replacing all omits by a new set of random responses chosen independently of the original set. No parameters were estimated from the CV data. Instead, the estimated $\log L$ was evaluated for the CV data, using estimated parameters in place of the unknown true values in (3).

For test M , when the estimated parameters were MLE's obtained from the original filled-in data, the estimated $\log L$ of (3) for the CV data was -52348 , approximately; when the estimated parameters were obtained from the raw data by (5), the estimated $\log L$ of (3) for the CV data was -52295 . Thus *the new method (5) in this case provides better estimates of the parameters of the conventional icc model (higher likelihood for the cross-validation data) than does the conventional method itself.*

The test V results support a similar conclusion. When the estimates were MLE's obtained from (3), $\log L$ in (3) for the CV data was -58277 ; when the estimates were obtained by (5), $\log L$ in (3) for the CV data was -58224 . Because of the time and expense involved, similar studies for tests v and m have not been made.

* Note added in proof: Samejima's work [1973] suggests that there might be more than one local maximum of the likelihood function for examinees such as these. A survey of the likelihood function for each of these 56 examinees disclosed no sign of additional local maxima.

The conclusions reached seem very plausible in view of the theoretical results proved in the Appendix. Further empirical confirmation would be desirable also.

7. *Artificial data.* Artificial data for 2995 hypothetical examinees on a hypothetical test composed of 90 five-choice items were generated so as to be very much like the real test V data, but with known item and person parameters. This was done by taking the parameters estimated from the test V data, modifying them slightly in an attempt to correct for errors of measurement, and then taking the modified parameters as known values from which artificial data were generated by use of the Hambleton-Rovinelli [1973] computer program.

Item and person parameters were then estimated from the data by means of (5). Because of the uncertainty already mentioned as to the ability level of very low-scoring examinees, the twelve values of θ_a^* below -4.0 were replaced by $\theta_a^* = -4.0$. Also, the a_i^* (estimated a_i) for 3 of the 90 items were held at $a_i^* = 1.75$ instead of allowing them to increase without limit.

In the case of easy items, the data usually provide no evidence at all as to the value of the c -parameter. In order to avoid absurd (although possibly harmless) estimates of c_i in such cases, all c_i^* (estimated c_i) were set equal to the median c_i^* of all items whenever the estimated standard error of the c_i^* was too large, as determined by a complicated set of rules [Wingersky & Lord, 1973] built into the computer program.

At the end of the study, the following correlations were found between the estimated and the true parameters. For 2995 examinees, $r_{\theta\theta} = .959$; for 90 items, $r_{b\theta} = .988$, $r_{a\theta} = .920$. Note that the squared correlation $r_{\theta\theta}^2 = .920$ may be compared with the Kuder-Richardson formula-20 reliability coefficient for number-right scores, which is .927.

Mean square errors of estimation would be meaningless here except in the case of the c -parameters. The reason is that the origin of the scale used to describe the θ_a and the b_i is completely arbitrary (it can be assigned at will without changing the numerical value of the likelihood function). The unit of measurement of the scales used to describe the θ_a , b_i , and a_i is equally arbitrary. For the 54 hardest items, the mean square error of the c_i^* is $(1/n) \sum (c_i^* - c_i)^2 = .0022$. The corresponding correlation is $r_{c\theta} = .851$. The lower accuracy of estimated c 's is neither surprising nor cause for concern, since for easy items the data will be equally well fitted regardless of c .

Appendix A

We are given answer sheets with three kinds of responses: rights, omits, and wrongs, which will be denoted by $v_i = 1, c$, and 0 , respectively. This appendix deals with the new estimation procedure for the case where there is just one examinee, the item parameters having been already determined.

Rewriting (4) with the subscript a dropped from most of the symbols, we have

$$(A1) \quad L_a^*(v_1, v_2, \dots, v_n \mid \theta) = \prod_{i=1}^n P_i^{v_i} Q_i^{1-v_i}.$$

Our estimate of the ability of the examinee tested is θ^* , the value of θ that maximizes (A1).

If we replace all omits by randomly assigned responses, the likelihood function for the resulting "filled-in" data under the icc model is

$$(A2) \quad L_a(u_1, u_2, \dots, u_n \mid \theta) = \prod_{i=1}^n P_i^{u_i} Q_i^{1-u_i},$$

where $u_i = 1$ or 0 . This equation is the same as (2) except that a has been dropped from most of the subscripts. The MLE $\hat{\theta}$ of the examinee's ability obtained from (A2) is justified by the icc model. We do not know how to compute MLE from the original data. The empirical results given in the last section suggest that θ^* computed from the original data may be a superior estimator to $\hat{\theta}$. It is most desirable, however, to have mathematical proofs of some of the properties of θ^* . No such proofs have been given so far. The purpose of this Appendix is to indicate some relationships between θ^* and $\hat{\theta}$ that hold for large n . It will be shown that θ^* (computed from the original data) is a consistent estimator with a sampling variance smaller than that of $\hat{\theta}$ (obtained after replacing omits by random responses).

In all that follows, we assume that θ is bounded, at least for any group of examinees that we consider testing. Another limitation is that we cannot estimate the ability of examinees who answer all items correctly or all items incorrectly. We deal only with examinees who give at least one right answer and at least one wrong answer. Since θ is bounded, the probability that an examinee will be excluded by this limitation approaches 0 for large n .

We assume that P_i is an increasing function of θ , twice differentiable, with $0 \leq c_i < P_i < 1$, c_i being the lower asymptote. These assumptions are easily satisfied by all icc ordinarily used for cognitive tests. We will avoid using extremely difficult or extremely easy items, so we can assume that P_i is bounded away from 0 and 1.

Let us reorder the items so that the s omitted items are numbered $i = 1, 2, \dots, s$. Then equations (A1) and (A2) become respectively

$$(A3) \quad L_a^* = \prod_{i=1}^s P_i^{c_i} Q_i^{1-c_i} \cdot \prod_{i=s+1}^n P_i^{u_i} Q_i^{1-u_i}$$

$$(A4) \quad L_a = \prod_{i=1}^s P_i^{u_i} Q_i^{1-u_i} \cdot \prod_{i=s+1}^n P_i^{u_i} Q_i^{1-u_i}.$$

Taking logarithms and dividing by n we can write

$$(A5) \quad \frac{1}{n} \log L_a^* = \frac{1}{n} \sum_{i=1}^s c \log \frac{P_i}{Q_i} + \frac{1}{n} \sum_{i=1}^n \log Q_i + \frac{1}{n} \sum_{i=s+1}^n u_i \log \frac{P_i}{Q_i},$$

$$(A6) \quad \frac{1}{n} \log L_a^* = \frac{1}{n} \log L_a - \frac{1}{n} \sum_{i=1}^s (u_i - c) \log \frac{P_i}{Q_i}.$$

Since the u_i are assigned at random with probability c that $u_i = 1$, the quantity Z defined by

$$(A7) \quad Z \equiv \frac{1}{s} \sum_{i=1}^s (u_i - c) \log \frac{P_i}{Q_i}$$

is the average of s observations, each on a random variable having a mean of zero. If $s \rightarrow \infty$, the variance of (A7) always $\rightarrow 0$ (since Q_i is bounded away from 0). Consequently the last term in (A6) always converges in probability to zero.

Thus, for large n the likelihood function (A2) and the new function (A1) tend to the same limit. This result makes the function (A1) a plausible function to investigate, even though (A1) is itself not a likelihood function.

Likelihood Equation

The log likelihood from (A2) is $\sum_i u_i \log P_i + \sum_i (1 - u_i) \log Q_i$. Taking the derivative with respect to θ gives the result

$$(A8) \quad \frac{d \log L_a}{d\theta} = \sum_{i=1}^n \left(\frac{u_i}{P_i} - \frac{1 - u_i}{Q_i} \right) P_i',$$

where $P_i' \equiv dP_i/d\theta$. Setting this derivative equal to zero yields a familiar likelihood equation

$$(A9) \quad \frac{d \log L_a}{d\theta} = \sum_{i=1}^n (u_i - P_i) \frac{P_i'}{P_i Q_i} = 0.$$

Similarly, setting $d \log L_a^*/d\theta = 0$, we obtain from (A1)

$$(A10) \quad \frac{d \log L_a^*}{d\theta} = \sum_{i=1}^n (v_i - P_i) \frac{P_i'}{P_i Q_i} = 0.$$

When the items are ordered with omitted items first, (A10) can be written

$$(A11) \quad \frac{d \log L_a^*}{d\theta} = \frac{d \log L_a}{d\theta} - \sum_{i=1}^s (u_i - c) \frac{P_i'}{P_i Q_i} = 0.$$

Consistency

Cramér's well-known proof [1946, section 33.3] that under regularity conditions a likelihood equation has a solution that converges in probability to the true value θ_0 as $n \rightarrow \infty$ applies with minor modifications to $\hat{\theta}$ obtained from (A9) (the cited proof only covers identically distributed variables). The same conclusion follows from Hoadley's [1971] theorem.

Cramér's approach can also be used to prove that θ^* , the solution of (A11), converges in probability to the true value θ_0 . When (A11) is divided by n and then $d \log L_a / d\theta$ expanded by a three-term Taylor's formula, we obtain after replacing θ by θ^*

$$(A12) \quad \frac{1}{n} \frac{d \log L_a^*}{d\theta} \Big|_{\theta^*} = B_0 - \frac{1}{n} \sum_{i=1}^s (u_i - c) \frac{P_i^{*'}}{P_i^* Q_i^*} + B_1(\theta^* - \theta_0) + \frac{1}{2} B_2(\theta^* - \theta_0)^2,$$

where $P_i^* \equiv P_i(\theta^*)$, et cetera, and

$$B_0 = \frac{1}{n} \frac{d \log L_a}{d\theta} \Big|_{\theta_0}, \quad B_1 = \frac{1}{n} \frac{d^2 \log L_a}{d\theta^2} \Big|_{\theta_0}, \quad B_2 = \frac{1}{n} \frac{d^3 \log L_a}{d\theta^3} \Big|_{\theta_1},$$

where θ_1 lies between θ^* and θ_0 . If, as we suppose, s/n approaches some constant less than 1 as $n \rightarrow \infty$, then the quantity $(1/n) \sum_{i=1}^s (u_i - c) P_i^{*'} / P_i^* Q_i^*$ can be combined with the first (lowest order) term and neglected when n is large. Cramér's argument then shows under regularity conditions that θ^* converges to θ_0 .

Sampling Variance and Efficiency

Let \mathcal{E} denote expectation over the population of items, so that $\mathcal{E}u_i = P_i^0$ and $\text{Var } u_i = P_i^0 Q_i^0$, where $P_i^0 = P_i(\theta_0)$, et cetera. The asymptotic sampling variance of the maximum likelihood estimator is

$$(A13) \quad \begin{aligned} \text{Var } \hat{\theta} &= \left[\mathcal{E} \left(\frac{d \log L_a}{d\theta} \right)^2 \Big|_{\theta=\theta_0} \right]^{-1} \\ &= \left[\mathcal{E} \left\{ \sum_{i=1}^n (u_i - P_i^0) \frac{P_i^{0'}}{P_i^0 Q_i^0} \right\}^2 \right]^{-1} \\ &= \left[\sum_{i=1}^n \sum_{j=1}^n \frac{P_i^{0'} P_j^{0'}}{P_i^0 Q_i^0 P_j^0 Q_j^0} \mathcal{E}(u_i - P_i^0)(u_j - P_j^0) \right]^{-1} \\ &= \left(\sum_{i=1}^n \frac{P_i^{0'2}}{P_i^{02} Q_i^{02}} \text{Var } u_i \right)^{-1} \\ &= \left(\sum_{i=1}^n \frac{P_i^{0'2}}{P_i^0 Q_i^0} \right)^{-1}. \end{aligned}$$

Thus, from (A12), following Cramér,

$$(A14) \quad \sqrt{n} (\theta^* - \theta_0) = \frac{\frac{1}{\sqrt{n}} \frac{d \log L_a}{d\theta} \Big|_{\theta_0} - \frac{1}{\sqrt{n}} \sum_{i=1}^s (u_i - c) \frac{P_i^{*'}}{P_i^* Q_i^*}}{-B_1 - \frac{1}{2} B_2 (\theta^* - \theta_0)}.$$

Our final problem is to find the asymptotic variance of $\sqrt{n} (\theta^* - \theta_0)$.

As n increases, $(\theta^* - \theta_0) \rightarrow 0$, so the second term in the denominator

converges in probability to zero. Thus, the entire denominator converges to

$$-\frac{1}{n} \varepsilon \left. \frac{d^2 \log L_a}{d\theta^2} \right|_{\theta_0} = \frac{1}{n \text{Var } \hat{\theta}} = k_0^2, \text{ say.}$$

Since $\text{Var } \hat{\theta}$ is of order $1/n$, k_0^2 is independent of n to our order of approximation. Since $s/n \rightarrow$ some constant as $n \rightarrow \infty$, the entire numerator is a random variable with zero mean and a finite variance, which we must now proceed to determine. The variance of the numerator will be obtained as the variance of the first term plus the variance of the second minus twice their covariance.

The term $(1/\sqrt{n}) (d \log L_a/d\theta) |_{\theta_0}$ is a random variable with zero mean and variance

$$(A15) \quad \frac{1}{n} \varepsilon \left(\frac{d \log L_a}{d\theta} \right)^2 \Big|_{\theta_0} = \frac{1}{n \text{Var } \hat{\theta}} = k_0^2.$$

The term $(1/\sqrt{n}) \sum_{i=1}^s (u_i - c)(P_i^{*'}/P_i^* Q_i^*)$ has zero mean. When s is fixed and large, its variance is approximately

$$(A16) \quad \frac{1}{n} \sum_{i=1}^s \varepsilon \left[(u_i - c)^2 \frac{P_i^{*'}^2}{P_i^{*2} Q_i^{*2}} \right] = \frac{1}{n} c(1 - c) \sum_{i=1}^s \frac{P_i^{0'2}}{P_i^{02} Q_i^{02}},$$

since the u_i , $i = 1, 2, \dots, s$ are not used in computing the P_i^* . For fixed s , the covariance between the two numerator terms in (A14) is found from (A9), using the same argument, to be

$$\begin{aligned} (A17) \quad & \frac{1}{n} \varepsilon \left[\sum_{i=1}^n (u_i - P_i^0) \frac{P_i^{0'}}{P_i^0 Q_i^0} \cdot \sum_{j=1}^s (u_j - c) \frac{P_j^{*'}}{P_j^* Q_j^*} \right] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^s \varepsilon \left[\{ (u_i - c) + (c - P_i^0) \} \frac{P_i^{0'}}{P_i^0 Q_i^0} (u_j - c) \frac{P_j^{*'}}{P_j^* Q_j^*} \right] \\ &= \frac{1}{n} \sum_{j=1}^s \frac{P_j^{0'}}{P_j^0 Q_j^0} \varepsilon \left[\frac{P_j^{*'}}{P_j^* Q_j^*} (u_j - c)^2 \right] \\ &= \frac{1}{n} c(1 - c) \sum_{j=1}^s \frac{P_j^{0'2}}{P_j^{02} Q_j^{02}}, \end{aligned}$$

approximately for large s and n .

We will need the general formulas for any random variables y, z, s :

$$(A18) \quad \text{Var}(y) = \varepsilon[\text{Var}(y | s)] + \text{Var}[\varepsilon(y | s)],$$

$$(A19) \quad \text{Cov}(y, z) = \varepsilon[\text{Cov}(y, z | s)] + \text{Cov}[\varepsilon(y | s), \varepsilon(z | s)].$$

Conveniently, the second term on the right of each formula is zero for the present application. Thus, by taking the expectation of (A16) with respect to s , we obtain the unconditional asymptotic variance

$$(A20) \quad \text{Var} \left[\frac{1}{\sqrt{n}} \sum_{i=1}^s (u_i - c) \frac{P_i^{*'}}{P_i^* Q_i^*} \right] = \frac{1}{n} c(1 - c) S,$$

where

$$(A21) \quad S \equiv \varepsilon \sum_{i=1}^s \frac{P_i^{0,2}}{P_i^{0,2} Q_i^{0,2}} = \sum_{i=1}^n \frac{P_i^{0,2} \omega_i}{P_i^{0,2} Q_i^{0,2}},$$

a positive quantity, where ω_i is the probability that the examinee will omit an item with characteristic curve P_i^0 . Similarly, the unconditional covariance corresponding to (A17) is found to be the same quantity, $(1/n)c(1 - c)S$.

From (A15) and (A20), we find the unconditional asymptotic variance of the numerator in (A14) to be

$$k_0^2 - \frac{1}{n} c(1 - c)S.$$

Finally, since the denominator is asymptotically a constant, k_0^2 , the asymptotic variance of $\sqrt{n}(\theta^* - \theta^0)$ is the variance of the numerator divided by k_0^4 or

$$(A22) \quad \text{Var } \theta^* = \text{Var } \hat{\theta}[1 - c(1 - c)S \text{ Var } \hat{\theta}] \\ = \frac{1 - c(1 - c)S / \left(\sum_{i=1}^n P_i^{0,2} / P_i^0 Q_i^0 \right)}{\sum_{i=1}^n P_i^{0,2} / P_i^0 Q_i^0}.$$

It is thus seen that the new method applied to the raw data ((4) and A1) has a smaller asymptotic sampling error than does the maximum likelihood method applied (of necessity) to the filled-in data ((2) and A2). The relative efficiency of the MLE is given by the term in brackets on the right of (A22).*

REFERENCES

- Birnbaum, A. Some latent trait models and their use in inferring an examinee's ability. In F. M. Lord & M. R. Novick, *Statistical theories of mental test scores*. Reading, Mass.: Addison-Wesley, 1968. Chapters 17-20.
- Bock, R. D. Estimating item parameters and latent ability when responses are scored in two or more nominal categories. *Psychometrika*, 1972, **37**, 29-51.
- Bock, R. D. & Lieberman, M. Fitting a response model for n dichotomously scored items. *Psychometrika*, 1970, **35**, 179-197.
- Cramér, H. *Mathematical methods of statistics*. Princeton, N. J.: Princeton University Press, 1946.
- Hambleton, R. K. & Rovinelli, R. A FORTRAN IV program for generating examinee response data from logistic test models. *Behavioral Science*, 1973, **18**, 74.
- Hoadley, B. Asymptotic properties of maximum likelihood estimators for the independent not identically distributed case. *The Annals of Mathematical Statistics*, 1971, **42**, 1977-1991.
- Kendall, M. G. & Stuart, A. *The advanced theory of statistics*. Vol. 1. New York: Hafner, 1958.

* The writer is indebted to Prof. Robert Jennrich for finding an error in an earlier version of this conclusion.

- Lord, F. M. An application of confidence intervals and of maximum likelihood to the estimation of an examinee's ability. *Psychometrika*, 1953, **18**, 57-76.
- Lord, F. M. An analysis of the Verbal Scholastic Aptitude Test using Birnbaum's three-parameter logistic model. *Educational and Psychological Measurement*, 1968, **28**, 989-1020.
- Lord, F. M. Estimating true-score distributions in psychological testing (An empirical Bayes estimation problem). *Psychometrika*, 1969, **34**, 259-299.
- Lord, F. M. Item characteristic curves estimated without knowledge of their mathematical form—a confrontation of Birnbaum's logistic model. *Psychometrika*, 1970, **35**, 43-50.
- Lord, F. M. Power scores estimated by item characteristic curves. *Educational and Psychological Measurement*, 1973, **33**, 219-224.
- Meredith, W. Some results based on a general stochastic model for mental tests. *Psychometrika*, 1965, **30**, 419-440.
- Samejima, F. A comment on Birnbaum's three-parameter logistic model in the latent trait theory. *Psychometrika*, 1973, **38**, 221-233.
- Slakter, M. J. The effect of guessing strategy on objective test scores. *Journal of Educational Measurement*, 1968, **5**, 217-221.
- Torgerson, W. S. *Theory and methods of scaling*. New York: Wiley, 1958.
- Wingersky, M. S.; Lees, D. M.; Lennon, V. & Lord, F. M. A computer program for estimating true-score distributions and graduating observed-score distributions. *Educational and Psychological Measurement*, 1969, **29**, 689-692.
- Wingersky, M. S. & Lord, F. M. A computer program for estimating examinee ability and item characteristic curve parameters when there are omitted responses. Research Memorandum 73-2. Princeton, N. J.: Educational Testing Service, 1973.

Manuscript received 6/28/73

Revised manuscript received 1/23/74