

DOCUMENT RESUME

ED 238 914

TM 870 740

AUTHOR Mislav, Robert J.
 TITLE Exploiting Auxiliary Information about Items in the Estimation of Rasch Item Difficulty Parameters.
 INSTITUTION Educational Testing Service, Princeton, N.J.
 SPONS AGENCY Office of Naval Research, Washington, D.C. Psychological Sciences Div.
 REPORT NO ETS-RR-87-26-ONR
 PUB DATE Jul 87
 CONTRACT N00014-85-K-0683
 NOTE 5lp.
 PUB TYPE Reports - Research/Technical (143)

EDRS PRICE MF01/PC03 Plus Postage.
 DESCRIPTORS *Bayesian Statistics; Difficulty Level; Estimation (Mathematics); Intermediate Grades; *Item Analysis; *Latent Trait Theory; *Mathematical Models; *Maximum Likelihood Statistics; Predictive Measurement; Regression (Statistics); Test Items

IDENTIFIERS California Achievement Tests; *Item Parameters; *Linear Logistic Test Model; Linear Models; Rasch Model

ABSTRACT

Standard procedures for estimating item parameters in Item Response Theory models make no use of auxiliary information about test items, such as their format or content, or the skills they require for solution. This paper describes a framework for exploiting this information, thereby enhancing the precision and stability of item parameter estimates and providing diagnostic information about items' operating characteristics. In the proposed model, final item parameter estimates represent a compromise between Linear Logistic Test Model estimates, where items with identical features would have identical estimates, and unrestricted maximum likelihood estimates. The principles were illustrated in a context for which a relatively simple approximation is available: empirical Bayes (EB) estimation of Rasch item difficulty parameters. Computation proceeded in three steps (1) unrestricted maximum likelihood estimates of item parameters; (2) point estimates of the regression parameters; and (3) final estimates of item parameters. A numerical example applied EB estimation procedures to the responses from 150 sixth graders on the Fractions subtest of the California Achievement Test. Three models, varying in their assumptions of item exchangeability, were fitted to the data. Analysis showed that auxiliary information about item features contributed as much information about item parameters as the likelihood function did. (Author/LPG)

 * Reproductions supplied by EDRS are the best that can be made *
 * from the original document. *

ED288914

EXPLOITING AUXILIARY INFORMATION ABOUT ITEMS IN THE ESTIMATION OF RASCH ITEM DIFFICULTY PARAMETERS

Robert J. Mislevy

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

- This document has been reproduced as received from the person or organization originating it.
- Minor changes have been made to improve reproduction quality.

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

This research was sponsored in part by the
Personnel and Training Research Programs
Psychological Sciences Division
Office of Naval Research, under
Contract No. N00014-85-K-0683

Contract Authority Identification Number
NR No. 150-539

Robert J. Mislevy, Principal Investigator



Educational Testing Service
Princeton, New Jersey

July 1987

Reproduction in whole or in part is permitted for
any purpose of the United States Government.

Approved for public release; distribution
unlimited.

TM 870 740

REPORT DOCUMENTATION PAGE

Form Approved
OMB No 0704-0188

| | | | |
|--|--|---|---------------------------------------|
| 1a. REPORT SECURITY CLASSIFICATION Unclassified | | 1b. RESTRICTIVE MARKINGS | |
| 2a. SECURITY CLASSIFICATION AUTHORITY | | 3. DISTRIBUTION / AVAILABILITY OF REPORT Approved for public release; distribution unlimited. | |
| 2b. DECLASSIFICATION / DOWNGRADING SCHEDULE | | | |
| 4. PERFORMING ORGANIZATION REPORT NUMBER(S) RR-87-26-ONR | | 5. MONITORING ORGANIZATION REPORT NUMBER(S) | |
| 6a. NAME OF PERFORMING ORGANIZATION Educational Testing Service | 6b. OFFICE SYMBOL (If applicable) | 7a. NAME OF MONITORING ORGANIZATION Personnel & Training Research Programs, Office of Naval Research (Code 1142PT), 90 N. Quincy Street | |
| 6c. ADDRESS (City, State, and ZIP Code) Princeton, NJ 08541 | | 7b. ADDRESS (City, State, and ZIP Code) Arlington, VA 22217-5000 | |
| 8a. NAME OF FUNDING / SPONSORING ORGANIZATION | 8b. OFFICE SYMBOL (If applicable) | 9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER N00014-85-K-0683 | |
| 8c. ADDRESS (City, State, and ZIP Code) | | 10. SOURCE OF FUNDING NUMBERS | |
| | | PROGRAM ELEMENT NO 61153N | PROJECT NO RR04204 |
| | | TASK NO RR04204-01 | WORK UNIT ACCESSION NO. NR 150-539 |
| 11. TITLE (Include Security Classification) Exploiting Auxiliary Information about Items in the Estimation of Rasch Item Difficulty Parameters (Unclassified) | | | |
| 12. PERSONAL AUTHOR(S) Robert J. Mislevy | | | |
| 13a. TYPE OF REPORT Technical | 13b. TIME COVERED FROM _____ TO _____ | 14. DATE OF REPORT (Year, Month, Day) June 1987 | 15. PAGE COUNT 47 |
| 16. SUPPLEMENTARY NOTATION | | | |
| 17. COSATI CODES | | 18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number) | |
| FIELD | GROUP | Empirical Bayes Exchangeability | |
| 05 | 09 | Collateral information Item response theory | |
| | | Hierarchical models Linear logistic test model | |
| 19. ABSTRACT (Continue on reverse if necessary and identify by block number) Standard procedures for estimating the item parameters in IRT models make no use of auxiliary information about test items, such as their format or content, or the skills they require for solution. This paper describes a framework for exploiting this information about items' operating characteristics. The principles are illustrated in a context for which a relatively simple approximation is available: empirical Bayes estimation of Rasch item difficulty parameters. | | | |
| 20. DISTRIBUTION / AVAILABILITY OF ABSTRACT <input type="checkbox"/> UNCLASSIFIED/UNLIMITED <input checked="" type="checkbox"/> SAME AS RPT <input type="checkbox"/> DTIC USERS | | 21. ABSTRACT SECURITY CLASSIFICATION Unclassified | |
| 22a. NAME OF RESPONSIBLE INDIVIDUAL Dr. Charles Davis | | 22b. TELEPHONE (Include Area Code) 202-696-4046 | 22c. OFFICE SYMBOL ONR 1142PT |

EXPLOITING AUXILIARY INFORMATION ABOUT ITEMS IN THE
ESTIMATION OF RASCH ITEM DIFFICULTY PARAMETERS

Robert J. Mislevy

This research was sponsored in part by the
Personnel and Training Research Programs
Psychological Sciences Division
Office of Naval Research, under
Contract No. N00014-85-K-0683

Contract Authority Identification Number
NR No. 150-539

Robert J. Mislevy, Principal Investigator

Educational Testing Service
Princeton, New Jersey

June 1987

Reproduction in whole or in part is permitted
for any purpose of the United States Government.

Approved for public release; distribution
unlimited.

Copyright © 1987. Educational Testing Service. All rights reserved.

Abstract

Standard procedures for estimating the item parameters in IRT models make no use of auxiliary information about test items, such as their format or content, or the skills they require for solution. This paper describes a framework for exploiting this information, thereby enhancing the precision and stability of item parameter estimates and providing diagnostic information about items' operating characteristics. The principles are illustrated in a context for which a relatively simple approximation is available: empirical Bayes estimation of Rasch item difficulty parameters.

Keywords: Empirical Bayes
Collateral information
Hierarchical models
Exchangeability
Item response theory
Linear logistic test model

Exploiting Auxiliary Information about Items in the
Estimation of Rasch Item Difficulty Parameters

Two active lines of research item in response theory (IRT) incorporate additional information into the process of parameter estimation, augmenting that conveyed by item responses alone. One line, motivated by statistical considerations, uses Bayesian procedures to obtain more accurate estimates of item and examinee parameters. Enhanced stability and lower mean squared errors can be achieved by assuming exchangeability over item parameters of a given type (e.g., difficulty parameters), effectively shrinking estimates toward their mean in inverse proportion to the degree of information available directly about them (Mislevy, 1986; Swaminathan & Gifford, 1982, 1985). A second line, motivated by psychological considerations, incorporates theories about specific skills or subtasks required to answer an item correctly. Scheiblechner (1972) and Fischer's (1973) Linear Logistic Test Model (LLTM) is a prime example; Rasch model item difficulty parameters are cast as linear combinations of more basic parameters that reflect the contributions of psychologically salient features of each item.

The purpose of this paper is to bring out a confluence of these two lines of research. The idea is to embed the LLTM in a

Bayesian framework, maintaining the notion that item features may indeed tell us something about item parameters, but admitting they may not tell us everything. Final item parameter estimates are a compromise between LLTM estimates, where items with identical features would have identical estimates, and unrestricted maximum likelihood estimates.

In order to focus on concepts rather than numerical procedures, we concentrate on a context for which a relatively simple approximation is available. The Rasch IRT model for dichotomous items is assumed; a linear regression model with normal, homoscedastic residuals is posited for item parameters given their salient features; and, with what is commonly called an empirical Bayes approximation, final item parameter estimates are calculated with maximum likelihood estimates of the regression model treated as known. The result is a simplified version of Smith's (1973) linear model with response-surface prior distributions.

The procedures are illustrated with data from a fractions test for junior high school students. Precision gains and diagnostic uses of the approach are discussed.

Background

This section briefly reviews the three components of an IRT model that incorporates auxiliary information about items. First is the item response model--specifically, in this presentation, the Rasch model. Following that are overviews of Bayesian estimation of item parameters and of the linear logistic test model.

The Rasch Model

Let x_{ij} denote the response of examinee i to item j , taking the value 1 if correct and 0 if not. The Rasch model (Rasch, 1960/1980) gives the probability of a correct response as

$$P_j(\theta_i) = P(x_{ij} = 1 | \theta_i, \beta_j) \\ = \exp(\theta_i - \beta_j) / [1 + \exp(\theta_i - \beta_j)] \quad , \quad (1)$$

where β_j characterizes the difficulty of item j and θ_i characterizes the ability of examinee i . Under the usual assumption of local independence, the probability of a vector pattern $\underline{x}_i = (x_{i1}, \dots, x_{in})'$ of responses to n items is

$$P(\underline{x}_i | \theta_i, \underline{\beta}) = \prod_j P_j(\theta_i)^{x_{ij}} Q_j(\theta_i)^{1-x_{ij}}, \quad (2)$$

where $Q_j(\theta) = 1 - P_j(\theta)$ and $\underline{\beta} = (\beta_1, \dots, \beta_n)'$. Assuming the independence of responses over examinees, the probability of the data matrix $\underline{X} = (\underline{x}_1, \dots, \underline{x}_N)'$ of N examinees is the product of expressions like Equation 2:

$$P(\underline{X} | \underline{\theta}, \underline{\beta}) = \prod_i P(\underline{x}_i | \theta_i, \underline{\beta}) \quad (3)$$

Once \underline{X} has been observed, Equation 3 is interpreted as a likelihood function, and provides a basis for estimating parameters. The literature offers a number of alternative procedures for doing so, including

- o joint maximum likelihood (JML), which finds values of $\underline{\beta}$ and each θ that, taken together, maximize Equation 3 (Wright & Panchapakesan, 1969);
- o conditional maximum likelihood (CML), which finds the maximizing value of $\underline{\beta}$ given examinees' total scores (Andersen, 1973); and

o marginal maximum likelihood (MML), which finds the maximizing value of $\underline{\beta}$ after integrating over a distribution of examinee parameters (Bock & Aitkin, 1981; Thissen, 1982).

These solutions provide similar estimates of $\underline{\beta}$ when neither the number of items or examinees is small; under appropriate assumptions they are asymptotically equivalent, consistent, and multivariate normal (for details see Haberman, 1977, on CML and JML, and De Leeuw & Verhelst, 1986, on CML and MML.)

We will have use for the normal approximation to MML in a subsequent section. The MML likelihood function is obtained from Equation 3 by marginalizing over the examinee distribution:

$$L_M(\underline{\beta}|\underline{X}) = \prod_i \int P(\underline{x}_i|\theta, \underline{\beta}) p(\theta) d\theta \quad , \quad (4)$$

where $p(\theta)$, the density function for examinee parameters, may be specified a priori (as in Bock and Aitkin, 1981, and Thissen, 1982) or estimated from the data (as in Cressie and Holland, 1983). When both the numbers of items and examinees are large, the likelihood function is approximately a product over items of independent normal distributions:

$$L_M(\underline{\beta}|\underline{X}) \propto \prod_j \exp[-(\beta_j - \hat{\beta}_j)^2 / 2\hat{\sigma}_j^2] \quad , \quad (5)$$

where $\hat{\beta}$ are MML estimates and $\hat{\sigma}_j$ are their estimated standard errors. (Large N is sufficient for multivariate normality, but large n is also necessary for independence.)

Bayesian Estimation

The simultaneous estimation of many parameters can often be improved when it is reasonable to consider subsets of parameters as exchangeable members of corresponding populations (Efron & Morris, 1975; Lindley & Smith, 1972). The subjective notion that parameters are "in some sense similar" implies a correlational structure on prior beliefs, which can be formalized by modeling the parameters as if they were a random sample from a population whose parameters are themselves imperfectly known. Data related directly to each individual parameter also conveys information about the higher-level population parameters; the population structure in turn provides information about the individual parameters.

In typical applications, resulting estimates of individual parameters are drawn toward the center of their distribution in inverse proportion to the amount of information available about

them directly. An intuitive justification of shrinkage is that unrestricted ML estimates contain sampling errors, so we would expect that the more extreme estimates reflect in part large sampling errors in that direction. This reasoning is consistent with the fact that the expected variance of ML estimates in such cases generally exceeds the variance of the true parameters.

Swaminathan and Gifford (1982) applied this idea to the Rasch model by assuming exchangeability over examinees and over items. In a Bayesian extension of JML, they provide estimation equations for the joint mode of $\underline{\beta}$ and $\underline{\theta}$ in the posterior distribution

$$p(\underline{\theta}, \underline{\beta} | \underline{X}) \propto P(\underline{X} | \underline{\theta}, \underline{\beta}) p(\underline{\theta}) p(\underline{\beta}) \quad , \quad (6)$$

where $p(\underline{\theta})$ and $p(\underline{\beta})$ are marginalizations over respective normal distributions, the parameters of which are estimated in part from the data. As expected, Swaminathan and Gifford's simulations showed the Bayesian estimates to be closer to their overall mean than unrestricted maximum likelihood estimates, and to have smaller mean squared error.

A similar extension of MML is described in Mislevy (1986). Marginalizing over $\underline{\theta}$ but not over the mean μ and standard

deviation ϕ of identical normal priors for the β 's, he gives estimation equations for the joint mode of $\underline{\beta}$, μ , and ϕ^2 in the posterior distribution

$$P(\underline{\beta}, \mu, \phi^2) \propto L_M(\underline{\beta} | \underline{X}) \times \prod_j p(\beta_j | \mu, \phi^2) \times p(\mu, \phi^2) \quad (7)$$

As with Swaminathan and Gifford's procedure, this approach also yields estimates of β 's that are closer to their estimated mean than those of the corresponding maximum likelihood procedure.

The Linear Logistic Test Model

In addition to positing a Rasch model for item responses, as in Equations 1 through 3, the LLTM assumes a linear model for the item parameters:

$$\beta_j = \sum_{k=1}^K q_{kj} \eta_k$$

$$= \underline{q}'_j \underline{\eta} \quad (8)$$

or, in matrix notation,

$$\underline{\beta} = \underline{Q}'\underline{\eta}$$

The basic parameters of the LLTM are η_k , $k = 1, \dots, K$. They reflect the additive contributions to item difficulty of selected item features. The vector q_j contains coefficients relating item j to basic parameters. In Fischer's (1973) calculus example, q indicated the number and the type of operations a pupil must carry out in order to solve a differentiation item. In Mitchell's (1983) analysis of Paragraph Comprehension subtests from the Armed Services Vocational Aptitude Battery, q conveyed semantic and lexicographic features of a question and an associated reading passage. The reader is referred to Fischer and Formann (1982) for additional applications of the LLTM.

Estimates of LLTM basic parameters can be obtained by suitable modification of JML, CML, and MML algorithms for the unconstrained Rasch model. Differences in $2 \log$ likelihood between the two models can be compared with the chi-square distribution on $n - K$ degrees of freedom, to test the significance of the constraints of the LLTM under the assumption that the unrestricted Rasch model is true.

Fischer and Formann (1982) note that the initial hope of explaining all reliable variation of item difficulties in terms of

basic parameters has not been fulfilled; rigorous tests of fit almost always reject the LLTM. This finding is consistent with what test developers have known for decades: two items written to test the same skill will differ in difficulty as a function of idiosyncratic features such as visual format and word choice.

Typically, however, a meaningful amount of variation can be explained. The proportion of variance of unconstrained estimates accounted for was 76 percent in Fischer's calculus test, and ranged from 66 to 96 percent in Mitchell's Paragraph Comprehension tests. Even though LLTM estimates $\hat{\beta} = Q'\hat{\eta}$ are not wholly acceptable as estimates of β , then, their ability to relate item performance to cognitive theory has proven useful in applications such as assessing treatment effects and modeling item bias. To the extent that LLTM does fit, it aids an understanding of just what makes items difficult. To the extent that it does not fit, departures indicate items that are unexpectedly hard or easy given the features that usually determine difficulty. Poor item construction or alternative response strategies can be detected in this way.

A Combined Model

Rationale

The assumption of exchangeability in the Bayesian estimation procedures described in a preceding section typically leads to item parameter estimates that are more stable and have lower mean squared errors. Strictly speaking, however, assuming exchangeability over all parameters of a given type, and consequently shrinking them all to the same center, is justified only if we have no prior information to distinguish among them. This is rarely the case in item parameter estimation. In vocabulary tests, for example, we know which words are frequently used and which ones are not; we expect the familiar words to be easier. In Fischer's calculus test, we would expect an item demanding several differentiation rules to be more difficult than one demanding only a subset of the same rules.

As Fischer and Formann (1982) point out, we cannot generally expect a few salient features to explain item parameters in toto. We can, however, express many of our prior beliefs in terms of such features. In particular, a model combining key aspects of the LLTM and the exchangeability concept of Bayesian estimation might consider as exchangeable only parameters of items with the same pedagogically or psychologically relevant features.

Shrinkage would then be observed toward the center of the subset to which an item belongs--as estimated from items of that type and possibly from other items as well, if they shared some features with it. This shrinkage could quite possibly be in the opposite direction from the center of the item set as a whole.

The General Form of the Model

Let the known (possibly vector-valued) quantity q_j represent auxiliary information about item j ; let $p(\beta|q)$ be the density function representing the distribution of β parameters for items with the same (generic) value of q . (The possibility that $p(\beta|q)$ may depend on unknown parameters is introduced below.) The posterior distribution of β , given the data \underline{X} and the auxiliary information $\underline{Q} = (q_1, \dots, q_n)$, is obtained as

$$p(\underline{\beta}|\underline{X}, \underline{Q}) \propto L_M(\underline{\beta}|\underline{X}) p(\underline{\beta}|\underline{Q}) \\ = \prod_i \int P(x_i | \theta, \underline{\beta}) p(\theta) d\theta \times \prod_j p(\beta_j | q_j) \quad . \quad (9)$$

An implementation of Equation 9 inspired by the LTM is to assume a linear regression model for $p(\beta|q)$ --a response-surface prior, as introduced by Smith (1973) in the context of linear models. With

\underline{Q} and $\underline{\eta}$ defined exactly as in the LLTM, we can approximate prior beliefs about item parameters as $MVN(\underline{Q}'\underline{\eta}, \phi^2 \underline{I})$. Considering $\underline{\eta}$ and ϕ^2 as additional unknown parameters, the marginal posterior is obtained as

$$P(\underline{\beta}, \underline{\eta}, \phi^2 | \underline{X}, \underline{Q}) \propto L_M \times \phi^{-m} \prod_j \exp[-(\beta_j - q_j' \underline{\eta})^2 / 2\phi^2] p(\underline{\eta}, \phi^2) \quad . \quad (10)$$

As in the LLTM, a linear model based on salient features gives the central tendency of items with the same features q_j , namely $\tilde{\beta}_j = q_j' \underline{\eta}$. Unlike the LLTM, however, variation of true parameters around these central values is anticipated.

Computational procedures for computing the posterior mode of $\underline{\beta}$, or of $\underline{\beta}$, μ , and ϕ^2 jointly, are readily obtained by generalizing the algorithms given in Mislevy (1986). The resulting solutions can be applied in the 2- and 3-parameter logistic models as well as for the Rasch model. The technical details of this solution are not central to the present paper, however; in order to focus upon concepts and applications, we now turn to a relatively simple computing approximation for the Rasch model.

A Computing Approximation for the Rasch Model

This section describes empirical Bayes (EB) estimation of Rasch item parameters, assuming normal linear regression on salient item features. Two simplifications are applied to the exact posterior distribution given in Equation 10. First, the marginal likelihood function of $\underline{\beta}$ is replaced by the normal approximation given in Equation 5. Second, MLE's of the population parameters $\underline{\eta}$ and ϕ^2 are treated as known, after they have been estimated from MLE's $\hat{\beta}_j$ with their standard errors $\hat{\sigma}_j$ treated as known. (It is this use of point estimates of population parameters that is commonly associated with the term "empirical Bayes.") The resulting approximation takes the following form:

$$\begin{aligned}
 p(\underline{\beta}|\underline{X},\underline{Q}) &\propto L_M(\underline{\beta}|\underline{X}) \times p(\underline{\beta}|\underline{Q}) \\
 &\propto L_M(\underline{\beta}|\underline{X}) \times \iint \prod_j p(\beta_j|q_j, \eta, \phi^2) p(\underline{\eta}, \phi^2) d\underline{\eta} d\phi^2 \\
 &\propto \prod_j \exp\left[\frac{-(\beta_j - \hat{\beta}_j)^2}{2\hat{\sigma}_j^2} \right] \times \prod_j \exp\left[\frac{-(\beta_j - q_j \hat{\eta})^2}{2\hat{\phi}^2} \right] .
 \end{aligned}$$

From this combination of a likelihood and prior that are both proportional to independent normal densities, independent normal posteriors follow (Box & Tiao, 1973, p. 74):

$$p(\underline{\beta} | \underline{X}, Q) \propto \prod_j \exp\left[\frac{-(\beta_j - \tilde{\beta}_j)^2}{2\tilde{\sigma}_j^2} \right],$$

where the means and variances are given by well-known formulas:

$$\tilde{\beta}_j = (\hat{\sigma}_j^{-2} \hat{\beta}_j + \hat{\phi}^{-2} q_j' \hat{\eta}) / (\hat{\sigma}_j^{-2} + \hat{\phi}^{-2}) \quad (11)$$

and

$$\tilde{\sigma}_j^2 = (\hat{\sigma}_j^{-2} + \hat{\phi}^{-2})^{-1} \quad (12)$$

Computation thus proceeds in three steps:

1. Unrestricted maximum likelihood estimates of item parameters
2. Point estimates of the regression parameters
3. Final estimates of item parameters

Step 1: Unrestricted maximum likelihood estimates of item parameters

Rasch item parameter estimates $\hat{\beta}_j$ and corresponding standard errors $\hat{\sigma}_j$ can be obtained with any of a number of widely-available computer programs. Numerical values and small-sample properties of JML, CML, and MML estimates certainly differ, but any suffice for our illustrative purposes. For long tests and many examinees, all support the approximation of the marginal likelihood as a product of independent normal distributions, with means given by maximum likelihood estimates and standard deviations given by the associated standard errors.

Step 2: Point estimates of the regression parameters

The regression structure for item parameters and the normal approximation for the marginal likelihood lead to the following system of regression equations:

$$\hat{\beta}_j = \beta_j + e_j \quad ,$$

where $(e_1, \dots, e_n) \sim \text{MVN}[\underline{0}, \text{diag}(\sigma_1^2, \dots, \sigma_n^2)]$, and

$$\beta_j = \underline{q}'_j \underline{\eta} + f_j \quad ,$$

where $(f_1, \dots, f_n) \sim \text{MVN}(\underline{0}, \phi^2 \underline{I})$. Taken together, they imply

$$\hat{\beta}_j = \underline{q}'_j \underline{\eta} + h_j \quad ,$$

where $(h_1, \dots, h_n) \sim \text{MVN}[\underline{0}, \text{diag}(\sigma_1^2 + \phi^2, \dots, \sigma_n^2 + \phi^2)]$.

MLE's for $\underline{\eta}$ and ϕ^2 can be obtained simultaneously by applying Dempster, Laird, and Rubin's (1977) EM algorithm. A special case of Braun and Jones' (1985) implementation was employed for the examples that appear in the following section. Using provisional estimates $\hat{\underline{\eta}}$ and $\hat{\phi}^2$, the E-step computes conditional expectations of the unknown item parameters:

$$\begin{aligned} \tilde{\beta}_j &= E(\beta_j | \hat{\beta}_j, \hat{\sigma}_j, \hat{\underline{\eta}}, \hat{\phi}^2) \\ &= (\hat{\phi}^{-2} \tilde{\beta}_j + \hat{\sigma}_j^{-2} \hat{\beta}_j) / (\hat{\phi}^{-2} + \hat{\sigma}_j^{-2}) \quad , \end{aligned}$$

where $\tilde{\beta}_j = \underline{q}'_j \hat{\underline{\eta}}$ is the (provisional) modeled mean for all items with same features as item j . The M-step uses these results to produce improved estimates:

$$\hat{\eta} = (Q'Q)^{-1}Q'\bar{\beta}_j$$

and

$$n\hat{\phi}^2 = \bar{\beta}_j'\bar{\beta}_j - \hat{\eta}'QQ'\hat{\eta}$$

Cycles of this type are repeated until convergence is attained.

Because the distribution of the hypothetical "complete data" $(\hat{\beta}, \hat{\beta})$, with parameters ϕ^2 and η , belongs to the exponential family if σ is assumed known, convergence to a unique maximum is assured (Dempster, Laird, & Rubin, 1977).

Step 3: Final estimates of item parameters

The posterior means and variances for the β 's that follow from our simplifying assumptions can be calculated as in Equations 11 and 12. The EB estimate $\bar{\beta}_j$ is thus a weighted average of the ML estimated $\hat{\beta}_j$ and the regression estimate $\bar{\beta}_j$. The relative weights are the precisions of the two estimates being combined, implying that ...

1. poorly-estimated $\hat{\beta}$'s shrink toward their predicted means more strongly than well-estimated $\hat{\beta}$'s;

2. if all $\hat{\beta}$'s are well-estimated in comparison with the estimated variation around their modeled means, little shrinkage occurs and $\tilde{\beta}$ approaches $\hat{\beta}$; and
3. if all $\hat{\beta}$'s are poorly-estimated in comparison with the expected variation around their modeled means, much shrinkage occurs and $\tilde{\beta}$ approaches $\bar{\beta}$.

Posterior precision, or $\tilde{\sigma}_j^{-2} = \hat{\sigma}_j^{-2} + \hat{\phi}^{-2}$, is the sum of precision about β_j conveyed directly through the likelihood function and that conveyed indirectly through knowledge about item features. By exploiting auxiliary information, then, the precision of item parameter estimates can be increased without to testing additional examinees.

Empirical Bayes estimates are distinguished most significantly from "true" Bayes estimates by their failure to account for uncertainty associated with η and ϕ^2 . The nature of the consequent differences is to overstate the apparent precision of the final EB item parameter estimates, while affecting their values only minimally. The posterior variances tend to be too small, and the distributions should be more platykurtic, like a t-distribution rather than the normal. The magnitude of these effects diminishes as η and ϕ^2 are better determined by the data. Larger N generally leads to greater

precision, but test length n and the matrix of cross-products $Q'Q$ are also important. These influences affect the precision of regression parameters and residual variance in much the same manner as in standard regression analyses.

A Numerical Example

This section applies EB estimation procedures to the 20-item Fractions subtest of the California Achievement Test (CAT), Level 3, Form A (Tiegs & Clark, 1970). The data are Rasch item difficulty estimates and standard errors, estimated from the responses of 150 sixth-grade students with the JML routine in Wright, Mead, and Bell's (1980) BICAL computer program. These values appear in Table 1, along with a specification of salient features of each item. These features, based on the CAT table of item specifications, are as follows:

1. Addition (ADD). The student must solve an addition problem involving one or more fractions and/or mixed numbers.
2. Subtraction (SUB). The student must solve a multiplication problem involving one or more fractions and/or mixed numbers.
3. Multiplication (MUL). The student must solve a multiplication problem involving one or more fractions and/or mixed numbers.
4. Division (DIV). The student must solve a division problem involving one or more fractions and/or mixed number.

5. Common denominators (CD). The student must find a common denominator for two fractions with unlike denominators.
6. Reduction (RED). The student must reduce a fraction or mixed number to lowest terms.

A sequence of three models was fit to these data:

Model 1: EB item parameter estimates were obtained under an assumption of global exchangeability. That is, all items were shrunk toward their common mean. The resulting estimates approximate the results of Swaminathan and Gifford's (1982) procedures.

Model 2: EB estimates were obtained under the assumption of exchangeability among items with the same features, based on Table 1.

Model 3: EB estimates were again obtained, after modifying the model along lines suggested by an examination of the estimates and residuals from Model 2.

 Insert Table 1 about here

Model 1: Twenty items, global exchangeability

Most applications of EB estimation involve shrinkage to the common center of the parameter set. This is accomplished in our

framework by using a vector of ones for Q . The results of such an analysis for the CAT Fractions test are presented in Table 2 and Figure 1. The grand mean toward which all estimates are shrunk is 0.00 (the result of the scaling convention used in BICAL); the estimated standard deviation $\hat{\phi}$ of the β 's with $\hat{\sigma}$ treated as $\underline{\sigma}$, is 1.71. This compares with a standard deviation of 1.74 for the $\hat{\beta}$'s, reflecting the expectation that a set of maximum likelihood estimates will be more dispersed than the set of parameters they estimate. Accordingly, under the assumption of exchangeability over all items, the EB estimates shrink toward their common mean.

 Insert Table 1 and Figure 1 about here

They do not shrink very much, though. If we define shrinkage for item j as $(\hat{\beta}_j - \bar{\beta}_j)/(\hat{\beta}_j - \underline{q}'\eta)$, then it is only about 2-percent on the average. The reason is that the estimated variance of β , about 2.92, is very large compared to the estimation error variance of the individual item parameters, about .06 on the average. Information from the likelihood function from a sample size of 150 is sufficient to overwhelm the information about interitem similarities, when the items are as dissimilar in difficulty as those in the Fractions test.

Model 2: Twenty items, exchangeability given salient features

A second model posits exchangeability for items with the same CAT specifications. The Q matrix in this case consisted of the columns of feature indicators given in Table 1. Estimates of η and ϕ are given in Table 3; item-level results are listed in Table 4 and illustrated in Figure 2.

 Insert Tables 3 and 4 and Figure 2 about here

The values of the regression parameters η shown in Table 3 are reasonably consistent with expectations. The values for addition, subtraction, multiplication, and division can be interpreted as values to which items exhibiting that feature only will be shrunk. Addition and subtraction show lower (easier) values than multiplication and division. The values for common denominators and fraction reduction are both positive, indicating additional difficulty for an item if this subskill is demanded in order to carry out the basic operation. The modeled mean for straight addition items, for example, is -2.75 ; the mean for addition items that also require reduction is $-2.75 + 1.90$, or $-.85$. Such addition items are nearly as hard as straight division items.

The residual standard deviation $\hat{\phi}$ under Model 2 is .58, much lower than the comparable value of 1.71 in Model 1 and closer to the typical standard error of about .3. EB item parameter estimates in Table 4 thus exhibit greater shrinkage--9 to 30 percent. Now that items within the smaller subsets over which exchangeability is assumed are in fact more similar, the structure contributes more information with which to improve item parameter estimates. Average posterior precision increases by roughly 25 percent, an amount equivalent to that attainable to testing about 40 more examinees.

Note that estimates now shrink toward the appropriate one of several predicted means rather than to a single overall mean. One item whose EB estimate moves away from the overall mean is item 8, the hardest of three straight subtraction items. Even though it was easier than average to begin with, the imposed exchangeability structure indicates that we would expect it to be easy based on the tasks it presents; in this particular data set, it may have been a bit harder than we might expect.

The last column in Table 4, labeled "standardized difference," gives the distance of an ML estimate from its predicted center, in standard deviation units:

$$\text{standardized difference} = \frac{\hat{\beta}_j - \bar{\beta}_j}{(\hat{\phi}^2 + \hat{\sigma}_j^2)^{1/2}}$$

By highlighting items that are unexpectedly far from their predicted means, these values can be useful for model modification. In conjunction with plots like Figure 2, they can reveal systematic departures from our expectations, which, upon reflection, lead us to modify the model.

Consider as an example the three straight subtraction items, 6, 7, and 8. As mentioned above, Item 8 is more difficult than modeled, to an extent that ranks it among the largest residuals in absolute value. The largest absolute residual, and in opposite direction, is the item in the same subset, namely item 7. This item is considerably easier than modeled. An inspection of item content offers an explanation: Item 7 asks for the solution of "1/6 - 1/6," which can be obtained without any knowledge of fractions at all. Despite its usefulness in ranking examinees, this item may not be tapping the skills the test is ostensibly attempting to measure. Further investigation reveals a similar phenomenon among straight division items, where Item 16 asks for the solution of "4/5 + 4/5." An atypically large negative

residual (easier than expected) for this item is balanced by an atypically large positive residual for another item (17) with the same features.

Further examination of items with large residuals reveals two items that are noticeably easier than expected for the same reason: while formally fractions items, both Item 1 (straight addition) and 6 (straight subtraction) require only a whole number operations with a fraction carried along. Failing to distinguish these items from straight addition or subtraction items that combine two actual fractions, Model 2 overpredicts the difficulty of Items 1 and 6.

A final anomaly appears in Figure 2, for Item 5. Item 5 is one of the harder items to begin with, but the regression model yields a higher-yet prediction, much higher than even the highest ML estimate observed. This is the only item requiring both the common denominator and reduction skills, and the higher prediction follows from the additivity of the model. The unappealing result suggests an interaction of sorts; while two additional subskills are required, it appears likely that examinees who possess the CD skill (the harder of the two) also possess the RED skill. Thus, incremental difficulty over straight addition when both are

present is not much over that expected from the common denominator subskill alone.

Model 3: Eighteen times, exchangeability given salient features

The final model illustrated here modified Model 2 in three ways:

1. Items 7 and 16, which could be solved by means of properties of operations alone, are eliminated from further consideration.
2. A column is added to the Q matrix reflecting a new salient feature: WN, or whole numbers only, applying to Items 1 and 6 which require just operations on whole numbers while a fraction is carried along.
3. To reflect the interaction of CD and RED observed for Item 5, its q value for RED has been changed from a 1 to a zero. That is, the difficulty parameters of addition items requiring CD and RED are now considered exchangeable with those of items requiring CD, the more difficult skill, alone.

The data for Model 3 are shown in Table 5. The results of the analysis are shown in Table 6 (regression parameter estimates), Table 7 (item-level results), and Figure 3 (a plot of ML, EB, and regression estimates). The revisions from model 2 reduced the residual standard deviation substantially, from .58 to

.23. This is about the same degree of precision as is available from the likelihood, so that EB estimates are roughly a 50-50 compromise between ML and regression estimates. Taking the approximate posterior variances at face value--recall that they are probably underestimated--we would conclude that the use of auxiliary information about items yields an increase in precision equivalent to doubling the size of the sample of examinees.

Insert Table 5, 6, and 7 about here

The average magnitude of standardized residuals is about the same as that from Model 2 because the denominator with which they are calculated decreased when the estimate of ϕ^2 decreased. Neither these residuals nor Figure 3 exhibit readily interpretable patterns of departures from the model.

Insert Figure 3 about here

As with any model-fitting procedure, the analysis that led to Model 3 capitalizes to some degree upon idiosyncratic features of the data at hand. Resulting estimates of precision are overly

optimistic for this reason in addition to the expedients employed by the estimation procedure. Any serious attempt to model item difficulties in the fractions domain would obviously require more data and more thought than were needed simply to illustrate computational procedures.

Discussion

The potential benefits of using auxiliary information about items in item parameter estimation are increased precision and diagnostic capabilities. In the numerical example in the preceding section, auxiliary information contributed as much information about item parameters as the likelihood function did. Conditional on the veracity of the assumed exchangeability structure, then, precision was increased by an amount equal to that attainable by doubling the number of examinees. Diagnostic checks revealed two items that might not be measuring the skills intended, by offering items that contained fractions but could be solved without manipulating them.

The plausibility of the exchangeability structure can also be verified with diagnostic checks. Two additional safeguards also mitigate the effects of specification errors at this stage. First, if the structure is badly in error and items assumed exchangeable turn out not be very similar, shrinkage will be

minimal (as in Model 1 of the example). Of course, minimal shrinkage does not necessarily signal misspecification or lack of exchangeability; all other things being equal, shrinkage decreases as N increases. Second, increasing the sample size of examinees leads to consistent item parameter estimates even if the exchangeability structure is flawed.

The simplified computing approximation used in this paper works best for the Rasch model, where it is needed least; even fairly small sizes give reasonably good item parameter estimates there. The same ideas can be applied more profitably to IRT models with more parameters, each less well-determined by data (e.g., the 3-parameter logistic model, and models for multiple-category item responses). The computational procedures for the general model are then required, since it may not be possible to obtain finite unrestricted ML estimates and their standard errors. No explicit averaging of ML and regression estimates can be accomplished in those cases, and Bayesian estimates must be obtained directly from item responses.

References

- Andersen, E. B. (1973). Conditional inference and models for measuring. Copenhagen: Danish Institute for Mental Health.
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: An application of an EM algorithm. Psychometrika, 46, 443-459.
- Box, G. E. P., & Tiao, G. C. (1973). Bayesian inference in statistical analysis. Reading, MA: Addison-Wesley.
- Braun, H. I., & Jones, D. H. (1985). Use of empirical Bayes methods in the study of the validity of academic predictors of graduate school performance. GRE Board Professional Report No. 79-13p and ETS Research Report 84-34. Princeton, NJ: Educational Testing Service.
- Cressie, N., & Holland, P.W. (1983). Characterizing the manifest probabilities of latent trait models. Psychometrika, 48, 129-141.
- de Leeuw, J., & Verhelst, N. (1986) Maximum likelihood estimation in generalized Rasch models. Journal of Educational Statistics, 11, 183-196.

- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm (with discussion). Journal of the Royal Statistical Society, Series B, 39, 1-38.
- Efron, B., & Morris, C. (1975). Data analysis using Stein's estimator and its generalizations. Journal of the American Statistical Association, 70, 311-319.
- Fischer, G. H. (1973). The linear logistic test model as an instrument in educational research. Acta Psychologica, 37, 359-374.
- Fischer, G. H., & Formann, A. K. (1982). Some applications of logistic latent trait models with linear constraints on the parameters. Applied Psychological Measurement, 6, 397-416.
- Haberman, S (1977). Maximum likelihood estimates in exponential response models. Annals of Statistics, 5, 815-841.
- Lindley, D. V., & Smith A. F. M. (1972). Bayes estimates for the linear model (with discussion). Journal of the Royal Statistical Society, Series B, 34, 1-41.
- Mislevy, R. J. (1986). Bayes modal estimation in item response models. Psychometrika, 51, 177-196.

- Mitchell, K. J. (1983). Cognitive processing determinants of item difficulty on the verbal subtests of the Armed Services Vocational Aptitude Battery. Technical Report 598. Alexandria, VA: U. S. Army Research Institute for the Behavioral and Social Sciences.
- Rasch, G. (1960/1980). Probabilistic models for some intelligence and attainment tests. Copenhagen: Danish Institute for Educational Research. Chicago: University of Chicago Press (reprint).
- Scheiblechner, H. (1972). Das lernen und losen komplexer denkaufgaben. Zeitschrift fur Experimentelle und Angewandte Psychologie, 19, 476-506.
- Smith, A. F. M. (1973). Bayes estimates in one-way and two-way models. Biometrika, 60, 319-329.
- Swaminathan, H., & Gifford, J. A. (1982). Bayesian estimation in the Rasch model. Journal of Educational Statistics, 7, 175-191.
- Swaminathan, H., & Gifford, J. A. (1985). Bayesian estimation in the two-parameter logistic model. Psychometrika, 50, 349-364.

Thissen, D. (1982). Marginal maximum likelihood estimation in the one-parameter logistic model. Psychometrika, 47, 175-186.

Tiegs, E., & Clark, W. (1970). The California Achievement Tests: 1970 edition. Monterey, CA: McGraw-Hill.

Wright, B. D., Mead, R. J., & Bell, S. R. (1980). BICAL: Calibrating items with the Rasch model. Research Memorandum 23C. Chicago: Statistical Laboratory, Department of Education, University of Chicago.

Wright, B. D., & Panchapekesan, N. (1969). A procedure for sample-free item analysis. Educational and Psychological Measurement, 29, 23-48.

Acknowledgments

This work was supported by Contract No. N00014-85-K-0683, project designation N^o 150-539, from Personnel and Training Research Programs, Psychological Sciences Division, Office of Naval Research. The author is grateful to Charles Lewis and Peter Pashley for their comments and suggestions, to Henry Braun and Bruce Kaplan for their assistance in applying the EM estimation procedure described in the example, and to Donna Lembeck and Maxine Kingston for the figures.

Table 1

Item Data and Salient Features: All Items

| Item | \hat{b} | $\hat{\sigma}$ | 1 ADD | 2 SUB | 3 MUL | 4 DIV | 5 CD | 6 RED |
|------|-----------|----------------|----------|----------|----------|----------|---------|----------|
| 1 | -3.73 | .31 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | -2.02 | .20 | 1 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1.45 | .28 | 1 | 0 | 0 | 0 | 1 | 0 |
| 4 | 1.16 | .26 | 1 | 0 | 0 | 0 | 1 | 0 |
| 5 | 1.63 | .31 | 1 | 0 | 0 | 0 | 1 | 1 |
| 6 | -2.42 | .21 | 0 | 1 | 0 | 0 | 0 | 0 |
| 7 | -3.23 | .27 | 0 | 1 | 0 | 0 | 0 | 0 |
| 8 | -1.05 | .18 | 0 | 1 | 0 | 0 | 0 | 0 |
| 9 | 1.28 | .27 | 0 | 1 | 0 | 0 | 1 | 0 |
| 10 | .30 | .21 | 0 | 1 | 0 | 0 | 0 | 1 |
| 11 | -.41 | .18 | 0 | 0 | 1 | 0 | 0 | 0 |
| 12 | -.80 | .18 | 0 | 0 | 1 | 0 | 0 | 0 |
| 13 | 2.22 | .38 | 0 | 0 | 1 | 0 | 0 | 1 |
| 14 | 1.72 | .31 | 0 | 0 | 1 | 0 | 0 | 1 |
| 15 | 1.41 | .28 | 0 | 0 | 1 | 0 | 0 | 1 |
| 16 | -1.35 | .18 | 0 | 0 | 0 | 1 | 0 | 0 |
| 17 | .26 | .21 | 0 | 0 | 0 | 1 | 0 | 0 |
| 18 | 1.28 | .27 | 0 | 0 | 0 | 1 | 0 | 1 |
| 19 | 1.41 | .28 | 0 | 0 | 0 | 1 | 0 | 1 |
| 20 | 1.05 | .25 | 0 | 0 | 0 | 1 | 0 | 1 |

Table 2
Item-Level Results from Model 1

| Items | $\hat{\beta}$ | $\hat{\sigma}$ | $\bar{\beta}$ | $\hat{\phi}$ | $\bar{\beta}$ | $\bar{\sigma}$ | Shrink- age | Standard- ized differ- ence |
|-------|---------------|----------------|---------------|--------------|---------------|----------------|----------------|--------------------------------------|
| 1 | -3.73 | 0.31 | 0.00 | 1.71 | -3.61 | 0.31 | 0.03 | -2.14 |
| 2 | -2.02 | 0.20 | 0.00 | 1.71 | -1.99 | 0.20 | 0.01 | -1.17 |
| 3 | 1.45 | 0.28 | 0.00 | 1.71 | 1.41 | 0.28 | 0.03 | 0.84 |
| 4 | 1.16 | 0.26 | 0.00 | 1.71 | 1.13 | 0.26 | 0.02 | 0.67 |
| 5 | 1.63 | 0.31 | 0.00 | 1.71 | 1.58 | 0.31 | 0.03 | 0.94 |
| 6 | -2.42 | 0.21 | 0.00 | 1.71 | -2.38 | 0.21 | 0.01 | -1.40 |
| 7 | -3.23 | 0.27 | 0.00 | 1.71 | -3.15 | 0.27 | 0.02 | -1.86 |
| 8 | -1.05 | 0.18 | 0.00 | 1.71 | -1.04 | 0.18 | 0.01 | -0.61 |
| 9 | 1.28 | 0.27 | 0.00 | 1.71 | 1.25 | 0.27 | 0.02 | 0.74 |
| 10 | 0.30 | 0.21 | 0.00 | 1.71 | 0.30 | 0.21 | 0.01 | 0.17 |
| 11 | -0.41 | 0.18 | 0.00 | 1.71 | -0.41 | 0.18 | 0.01 | -0.24 |
| 12 | -0.80 | 0.18 | 0.00 | 1.71 | -0.79 | 0.18 | 0.01 | -0.46 |
| 13 | 2.22 | 0.38 | 0.00 | 1.71 | 2.12 | 0.37 | 0.05 | 1.27 |
| 14 | 1.72 | 0.31 | 0.00 | 1.71 | 1.67 | 0.31 | 0.03 | 0.99 |
| 15 | 1.41 | 0.28 | 0.00 | 1.71 | 1.37 | 0.28 | 0.03 | 0.81 |
| 16 | -1.35 | 0.18 | 0.00 | 1.71 | -1.34 | 0.18 | 0.01 | -0.78 |
| 17 | 0.26 | 0.21 | 0.00 | 1.71 | 0.26 | 0.21 | 0.01 | 0.15 |
| 18 | 1.28 | 0.27 | 0.00 | 1.71 | 1.25 | 0.27 | 0.02 | 0.74 |
| 19 | 1.41 | 0.28 | 0.00 | 1.71 | 1.37 | 0.28 | 0.03 | 0.81 |
| 20 | 1.05 | 0.25 | 0.00 | 1.71 | 1.03 | 0.25 | 0.02 | 0.61 |

Table 3

Estimates of Regression Parameters under Model 2

| Effect (η) | Estimate |
|-------------------------------|----------|
| 1. Addition | -2.75 |
| 2. Subtraction | -2.08 |
| 3. Multiplication | -.34 |
| 4. Division | -.61 |
| 5. Common denominators | 3.50 |
| 6. Reduction | 1.90 |
| Standard deviation (ϕ) | .58 |

Table 4
Item-Level Results from Model 2

| Items | $\hat{\beta}$ | $\hat{\sigma}$ | $\tilde{\beta}$ | $\hat{\phi}$ | $\tilde{\beta}$ | $\tilde{\sigma}$ | Shrink- age | Standard- ized differ- ence |
|-------|---------------|----------------|-----------------|--------------|-----------------|------------------|----------------|--------------------------------------|
| 1 | -3.73 | 0.31 | -2.75 | 0.58 | -3.61 | 0.27 | 0.22 | -1.49 |
| 2 | -2.02 | 0.20 | -2.75 | 0.58 | -2.10 | 0.19 | 0.11 | 1.19 |
| 3 | 1.45 | 0.28 | 0.75 | 0.58 | 1.32 | 0.25 | 0.19 | 1.09 |
| 4 | 1.16 | 0.26 | 0.75 | 0.58 | 1.09 | 0.24 | 0.17 | 0.64 |
| 5 | 1.63 | 0.31 | 2.65 | 0.58 | 1.86 | 0.27 | 0.22 | -1.55 |
| 6 | -2.42 | 0.21 | -2.08 | 0.58 | -2.38 | 0.20 | 0.12 | -0.55 |
| 7 | -3.23 | 0.27 | -2.08 | 0.58 | -3.02 | 0.24 | 0.18 | -1.80 |
| 8 | -1.05 | 0.18 | -2.08 | 0.58 | -1.14 | 0.17 | 0.09 | 1.69 |
| 9 | 1.28 | 0.27 | 1.42 | 0.58 | 1.30 | 0.24 | 0.18 | -0.22 |
| 10 | 0.30 | 0.21 | -0.18 | 0.58 | 0.24 | 0.20 | 0.12 | 0.77 |
| 11 | -0.41 | 0.18 | -0.34 | 0.58 | -0.40 | 0.17 | 0.09 | -0.11 |
| 12 | -0.80 | 0.18 | -0.34 | 0.58 | -0.76 | 0.17 | 0.09 | -0.75 |
| 13 | 2.22 | 0.38 | 1.56 | 0.58 | 2.02 | 0.32 | 0.30 | 0.96 |
| 14 | 1.72 | 0.31 | 1.56 | 0.58 | 1.68 | 0.27 | 0.22 | 0.25 |
| 15 | 1.41 | 0.28 | 1.56 | 0.58 | 1.44 | 0.25 | 0.19 | -0.23 |
| 16 | -1.35 | 0.18 | -0.61 | 0.58 | -1.29 | 0.17 | 0.09 | -1.21 |
| 17 | 0.26 | 0.21 | -0.61 | 0.58 | 0.16 | 0.20 | 0.12 | 1.42 |
| 18 | 1.28 | 0.27 | 1.29 | 0.58 | 1.28 | 0.24 | 0.18 | -0.01 |
| 19 | 1.41 | 0.28 | 1.29 | 0.58 | 1.39 | 0.25 | 0.19 | 0.19 |
| 20 | 1.05 | 0.25 | 1.29 | 0.58 | 1.09 | 0.23 | 0.16 | -0.37 |

Table 5

Item Data and Salient Features: Reduced Set

| Item | \hat{b} | $\hat{\sigma}$ | 1 ADD | 2 SUB | 3 MUL | 4 DIV | 5 CD | 6 RED | 7 WN |
|------|-----------|----------------|----------|----------|----------|----------|---------|----------|---------|
| 1 | -3.73 | .31 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| 2 | -2.02 | .20 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 1.45 | .28 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| 4 | 1.16 | .26 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| 5 | 1.63 | .31 | 1 | 0 | 0 | 0 | 1 | 0 | 0 |
| 6 | -2.42 | .21 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| (7) | | | | | | | | | |
| 8 | -1.05 | .18 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 9 | 1.28 | .27 | 0 | 1 | 0 | 0 | 1 | 0 | 0 |
| 10 | .30 | .21 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 11 | -.41 | .18 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 12 | -.80 | .18 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 13 | 2.22 | .38 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 14 | 1.72 | .31 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 15 | 1.41 | .28 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| (16) | | | | | | | | | |
| 17 | .26 | .21 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 18 | 1.28 | .27 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| 19 | 1.41 | .28 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| 20 | 1.05 | .25 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |

Table 6

Estimates of Regression Parameters under Model 3

| <u>Effect (η)</u> | <u>Estimate</u> |
|-----------------------------------|-----------------|
| 1. Addition | -1.90 |
| 2. Subtraction | -1.28 |
| 3. Multiplication | -.32 |
| 4. Division | -.25 |
| 5. Common denominators | 3.10 |
| 6. Reduction | 1.71 |
| 7. Whole numbers only | -1.41 |
| Standard deviation (ϕ) | .23 |

Table 7
Item-Level Results from Model 3

| Items | $\hat{\beta}$ | $\hat{\sigma}$ | $\bar{\beta}$ | $\hat{\sigma}$ | $\bar{\beta}$ | $\bar{\sigma}$ | Shrink- age | Standard- ized Differ- ence |
|-------|---------------|----------------|---------------|----------------|---------------|----------------|----------------|--------------------------------------|
| 1 | -3.73 | 0.31 | -3.30 | 0.23 | -3.45 | 0.18 | 0.65 | -1.11 |
| 2 | -2.02 | 0.20 | -1.89 | 0.23 | -1.96 | 0.15 | 0.44 | -0.43 |
| 3 | 1.45 | 0.28 | 1.21 | 0.23 | 1.30 | 0.18 | 0.61 | 0.67 |
| 4 | 1.16 | 0.26 | 1.21 | 0.23 | 1.19 | 0.17 | 0.57 | -0.14 |
| 5 | 1.63 | 0.31 | 1.21 | 0.23 | 1.35 | 0.18 | 0.65 | 1.10 |
| 6 | -2.42 | 0.21 | -2.70 | 0.23 | -2.55 | 0.15 | 0.47 | 0.89 |
| (7) | | | | | | | | |
| 8 | -1.05 | 0.18 | -1.28 | 0.23 | -1.14 | 0.14 | 0.39 | 0.80 |
| 9 | 1.28 | 0.27 | 1.82 | 0.23 | 1.60 | 0.17 | 0.59 | -1.53 |
| 10 | 0.30 | 0.21 | -0.32 | 0.23 | 0.36 | 0.15 | 0.47 | -0.42 |
| 11 | -0.41 | 0.18 | -0.32 | 0.23 | -0.38 | 0.14 | 0.39 | -0.30 |
| 12 | -0.80 | 0.18 | 1.39 | 0.23 | -0.61 | 0.14 | 0.39 | -1.65 |
| 13 | 2.22 | 0.38 | 1.39 | 0.23 | 1.60 | 0.19 | 0.74 | 1.89 |
| 14 | 1.72 | 0.31 | 1.39 | 0.23 | 1.50 | 0.18 | 0.65 | 0.87 |
| 15 | 1.41 | 0.28 | 1.39 | 0.23 | 1.39 | 0.18 | 0.61 | 0.07 |
| (16) | | | | | | | | |
| 17 | 0.26 | 0.21 | -0.25 | 0.23 | 0.02 | 0.15 | 0.47 | 1.66 |
| 18 | 1.28 | 0.27 | 1.46 | 0.23 | 1.38 | 0.17 | 0.59 | -0.50 |
| 19 | 1.41 | 0.28 | 1.46 | 0.23 | 1.44 | 0.18 | 0.61 | -0.13 |
| 20 | 1.05 | 0.25 | 1.46 | 0.23 | 1.27 | 0.17 | 0.55 | -1.21 |

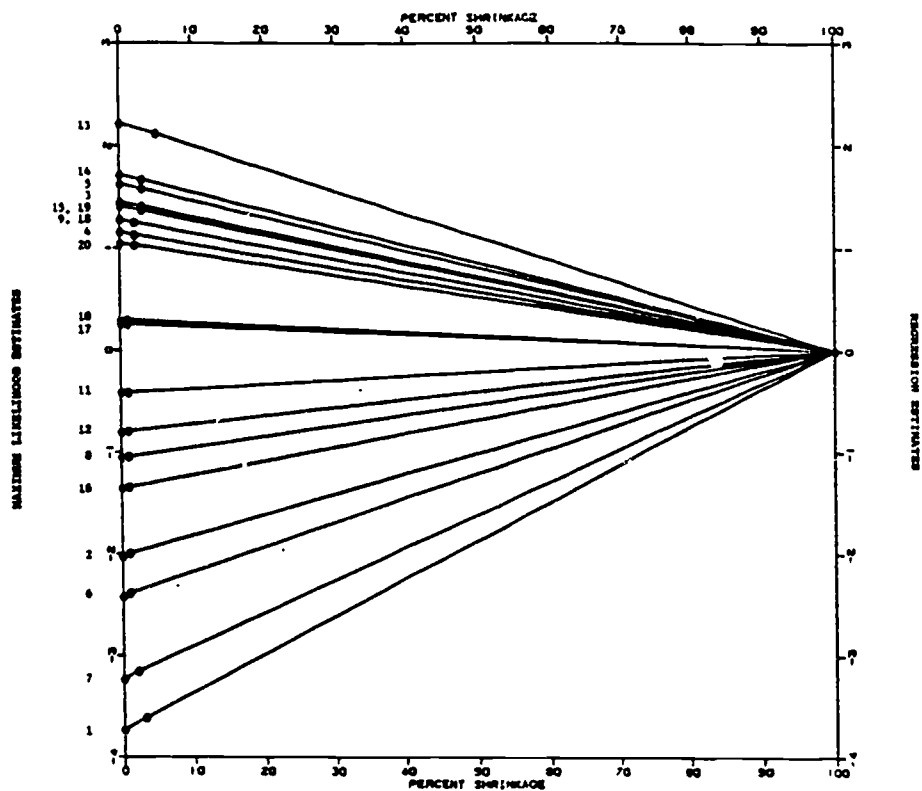


Figure 1: Maximum likelihood, regression, and empirical Bayes item parameter estimates: Model 1.

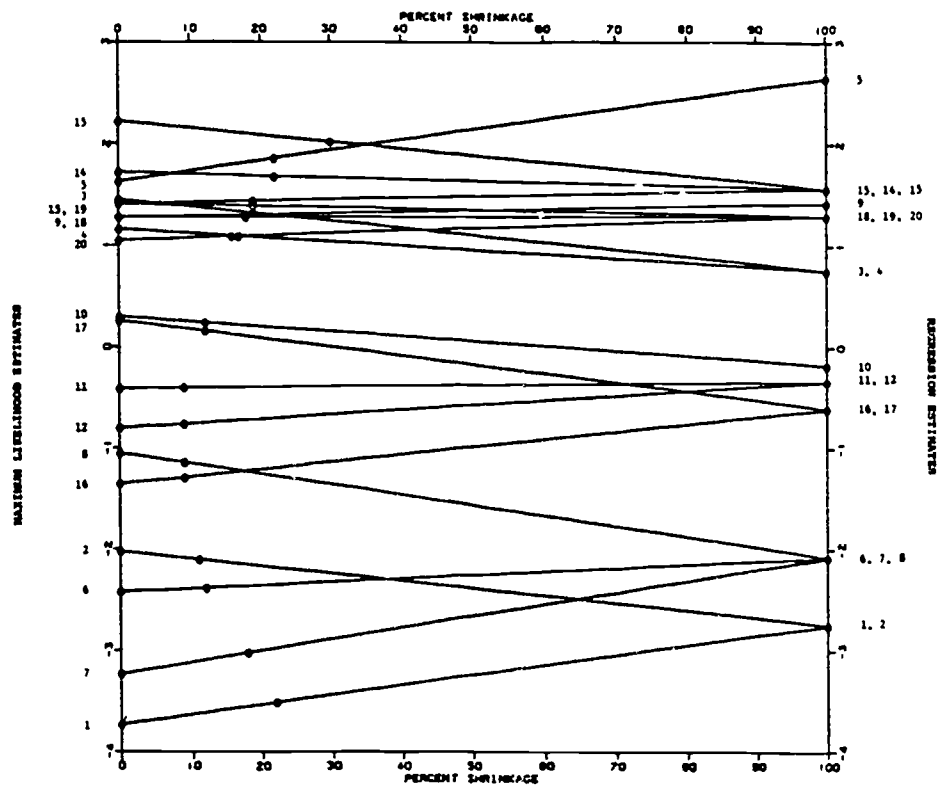


Figure 2: Maximum likelihood, regression, and empirical Bayes item parameter estimates: Model 2.

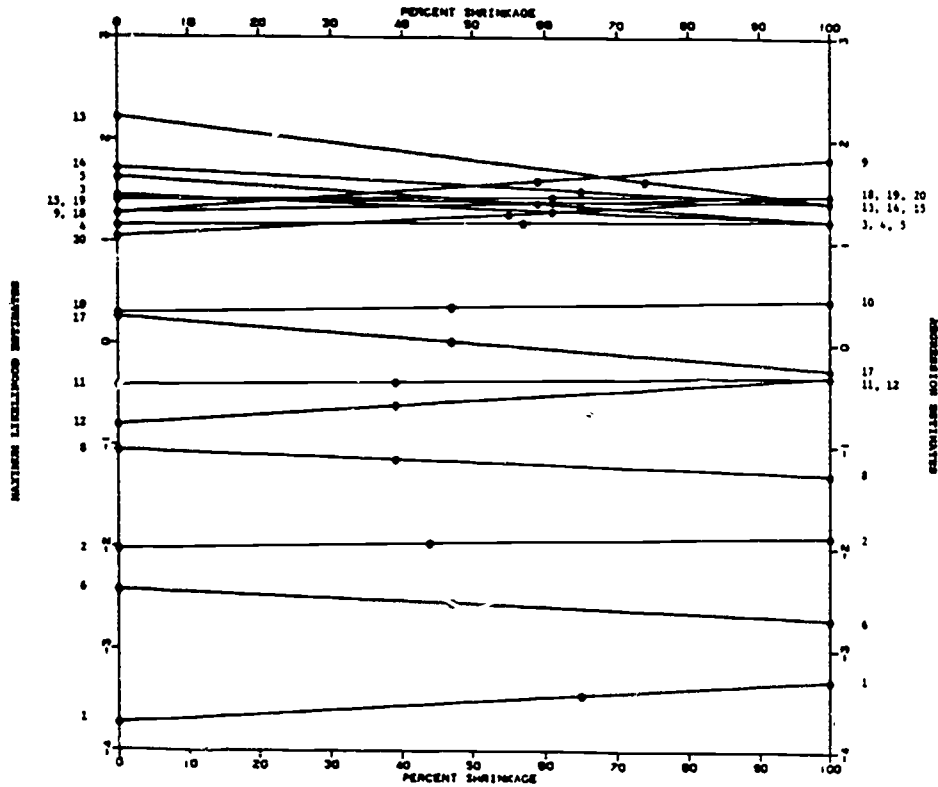


Figure 3. Maximum likelihood, regression, and empirical Bayes item parameter estimates: Model 3.