

# Inferring Examinee Ability When Some Item Responses Are Missing

Robert J. Mislevy<sup>1</sup> and Pao-Kuei Wu<sup>1,2</sup>

<sup>1</sup>Educational Testing Service

<sup>2</sup>University of California at Berkeley

The basic equations of item response theory (IRT) provide a foundation for inferring examinee's abilities and items' operating characteristics from observed responses. In practice, though, examinees will usually not have provided a response to every available item—for reasons that may or may not have been intended by the test administrator, and that may or may not be related to examinee ability. The mechanisms that produce missingness must be taken into account if correct inferences are to be drawn. Using concepts introduced by Rubin (1976), we discuss the implications for ability and item parameter estimation that are entailed by alternate test forms, targeted testing, adaptive testing, time limits, and omitted responses.

*Keywords:* Adaptive testing; Item response theory; Missing data; Omitted responses; Targeted testing

## Introduction

The capability to measure different examinees with different test items is an oft-cited advantage of item response theory (IRT). This option implies a problem of inference in the presence of missing data, since an examinee may not have provided a response to every item in the complete item set. Five types of missingness are in fact encountered regularly in routine applications of IRT:

Case 1: Alternate test forms. Two or more tests with similar content but different items are often employed to minimize carry-over effects (as in test-retest designs), reduce fatigue and practice effects (by splitting a test into shorter subtests), or avoid cheating behavior. A examinee is typically administered one from selected at random.

Case 2: Targeted testing. Two or more tests with similar content, but pitched at different levels of difficulty, can be used to make testing more efficient when background information (such as grade or courses taken) is available for deciding which test to administer to each examinee.

Case 3: Adaptive testing. Testing can also be made more efficient and less time-consuming if each item presented to an examinee is selected on the basis of his responses up to that point, and possibly background information as well.

Case 4: Not-reached items. Under typical testing conditions, some examinees will not reach the last few items on a test because of the time limit.

Case 5: Omitted items. Even when an item has been presented to an examinee and he has time to reach it, he will sometimes choose not to respond.

When incomplete data of any of these types are encountered, the IRT model that presumably accounts for the responses that are observed, is embedded in a more encompassing model that determines which responses will be observed and which will be missing. This paper discusses the implications that missing responses hold for likelihood and Bayesian inferences about examinee ability parameters and item parameters, assuming an IRT model holds. When can the process that causes missingness be ignored? When it cannot be ig-

---

This paper results from work carried out during the second author's summer internship at ETS. The first author was supported by Contract No. N00014-85-K-0683, project designation NR 150-539, from the Cognitive Science Program, Cognitive and Neural Sciences Division, Office of Naval Research. Reproduction in whole or in part is permitted for any purpose of the United States Government. We are grateful to Murray Aitkin and Kentaro Yamamoto for their comments and suggestions.

nored, how can it be modeled? How can conventional IRT methods for missing responses be evaluated in this framework?

The following section extends IRT notation to handle missingness, using concepts and notation from Little and Rubin (1987) and Rubin (1976). Next, Rubin's (1976) conditions for when the missingness process can be ignored are reviewed. Each of the five types of missingness listed above are then discussed in some detail in the problem of inferring ability when item parameters are known. This is followed by the extension to item parameter estimation. A final section summarizes our results.

### Background and Notation

At the heart of IRT is the model for the response to item  $j$ , with its possibly vector-valued parameter  $\beta_j$ , from an examinee with ability  $\theta$ . The Rasch model for dichotomous items, for example, posits

$$\mathbb{P}(U_j = u_j \mid \theta, b_j) = \frac{\exp\{u_j(\theta - b_j)\}}{1 + \exp\{\theta - b_j\}},$$

where  $u_j = 1$  denotes a correct response and  $u_j = 0$  an incorrect one, and  $b_j$  is the difficulty parameter of item  $j$ . We assume IRT functions that are twice differentiable, and interpret  $\mathbb{P}(U_j = 1 \mid \theta, \beta_j)$  as the proportion of correct responses we would expect to many items with  $\beta = \beta_j$  from many examinees with that value of  $\theta$ .

Under the usual assumption of local independence, the conditional probability of the response vector  $\mathbf{U} = (U_1, \dots, U_n)^\top$  for  $n$  items is obtained by the product rule:

$$\mathbb{P}(\mathbf{U} = \mathbf{u} \mid \theta, \beta) = \prod_{j=1}^n \mathbb{P}(U_j = u_j \mid \theta, \beta_j). \quad (1)$$

It is further assumed that if  $y$  denotes background information about an examinee such as age, GPA, or courses taken, then

$$\mathbb{P}(\mathbf{U} = \mathbf{u} \mid \theta, \beta, y) = \mathbb{P}(\mathbf{U} = \mathbf{u} \mid \theta, \beta).$$

When there is no possibility of missing responses, Equation (1) can be interpreted as a likelihood function, say  $L(\theta \mid \tilde{\mathbf{u}})$ , once a particular value  $\tilde{\mathbf{u}}$  of  $\mathbf{U}$  has been observed. *Direct likelihood inferences* are based solely on relative values of  $L$  at different values of  $\theta$ . It might be

said, for example, that the probability of  $\tilde{\mathbf{u}}$  is twice as high at  $\theta'$  than at  $\theta''$ . The maximum likelihood estimate (MLE),  $\hat{\theta}$ , is the value at which  $\tilde{\mathbf{u}}$  has the highest probability. Note that in direct likelihood inference, the MLE concerns only the data that were actually observed.

The role of the MLE in *sampling distribution inferences* concerns its distribution under repeated sampling of observations with a fixed "true" parameter value. If  $n$  is large, the sampling distribution of  $\hat{\theta}$  as computed from repeated observations of  $\mathbf{U}$  can be approximated by a normal distribution with mean  $\theta$  and variance

$$\sigma^2 = - \left[ \frac{\partial^2 \ell(\theta \mid \tilde{\mathbf{u}})}{\partial \theta \partial \theta^\top} \right]^{-1}$$

where  $\ell(\theta \mid \tilde{\mathbf{u}}) = \log L(\theta \mid \tilde{\mathbf{u}})$ . By considering the distribution of  $\hat{\theta}$  over hypothetical draws from the sample space, sampling distribution inferences involve datasets that could have been observed, but were not.

*Bayesian inferences* are based on the poosterior distribution for  $\theta$  given  $\tilde{\mathbf{u}}$ , or

$$p(\theta \mid \tilde{\mathbf{u}}) = k L(\theta \mid \tilde{\mathbf{u}}) p(\theta), \quad (2)$$

where  $k$  is a normalizing constant and  $p(\theta)$  conveys knowledge about  $\theta$  before a value of  $\mathbf{U}$  is observed. The posterior mean and mode of  $\theta$  are sometimes taken as point estimates in IRT. The posterior variance is approximated by  $\sigma^2$  when  $n$  is large. (This is the variance of the posterior distribution for  $\theta$  induced by the data actually observed, in contrast to the variance of an estimator over hypothetical repeated observations.)

In many applications of IRT, an examinee provides responses to only a subset of the  $n$  items to which responses could have been observed. The data thus consist of (i) the identification of the subset of items to which responses are observed and (ii) the responses to those items. The first inferential problem we address is to estimate an individual examinee's  $\theta$  from this extended observation, assuming that both the IRT model and the item parameters are known. To this end, we adapt notation from Little and Rubin (1987) and Rubin (1976) in defining the following terms:

- $\mathbf{u} = (u_1, \dots, u_n)^\top$  is the (hypothetical) random vector of responses to all items in the full item set.
- $\mathbf{m} = (m_1, \dots, m_n)^\top$  is an associated "missing data indicator" with each element taking values of 0 or 1. If  $m_j = 1$ , the value of  $U_j$  will be observed; if  $m_j = 0$ , the value of  $U_j$  will be missing.

•  $\mathbf{v} = (v_1, \dots, v_n)^\top$  conveys the data that are actually observed:  $v_j = u_j$  if  $m_j = 1$  but  $v_j = *$  if  $m_j = 0$ .

An observed value of  $\mathbf{m}$ , say  $\tilde{\mathbf{m}}$ , effects a partition of  $\mathbf{u}$ ,  $u$ ,  $\mathbf{v}$ , and  $v$  according to which elements are observed and which are missing. That is, we may write  $\mathbf{u} = (\mathbf{u}_{\text{mis}}, \mathbf{u}_{\text{obs}})$  to distinguish the missing and observed elements of  $\mathbf{u}$ , respectively. Similarly,  $u = (u_{\text{mis}}, u_{\text{obs}})$ ,  $\mathbf{v} = (\mathbf{v}_{\text{mis}}, \mathbf{v}_{\text{obs}})$ , and  $v = (v_{\text{mis}}, v_{\text{obs}})$ . As with  $\tilde{\mathbf{u}}$  and  $\tilde{\mathbf{m}}$ , let  $\tilde{v}$  denote a realized value of  $v$ .

### Example

An examinee is administered a two-item test. With each item scored right or wrong (1 or 0), there are  $2^2 = 4$  possible patterns for  $u$ : (0, 0), (0, 1), (1, 0), and (1, 1). The second response may be missing, however. With 1 representing “observed” and 0 representing “missing”, there are 4 conceivable patterns for  $m$ , of which (1, 0) and (1, 1) can be realized. If the examinee would have responded incorrectly to the first item and correctly to the second, but the response for the second item is missing, then  $\tilde{u} = (0, 1)$ ,  $\tilde{m} = (1, 0)$ , and  $\tilde{v} = (0, *)$ .

Inferences must of course be based on the data that are actually observed, namely realizations of  $\mathbf{v} = (j_{\text{obs}}, m)$ . Modeling the hypothetical complete data vector  $(u, m)$ —even if there is no intention of observing a response to every item—is a convenient way to begin. It forces us to explicate our beliefs about the relationships among ability, item response, and missingness—exactly what is required for building a sensible model for  $v$ . Recalling that  $p(u, m)$  can be written as  $p(m | u)p(u)$  or as  $p(u | m)p(m)$ , define the following densities:

•  $f_\theta(u)$  is the density for all  $n$  responses. In this paper,  $f_\theta(u)$  takes the form shown in Equation (1), so by local independence we can write  $f_\theta(u) = f_\theta(u')f_\theta(u'')$  for any ordering and partitioning of the items into  $(u', u'')$ —including  $(u_{\text{mis}}, u_{\text{obs}})$ .

•  $g_\phi(m | u)$  is the probability that  $\mathbf{m}$  takes value  $\mathbf{m} = (m_1, \dots, m_n)^\top$  given that  $\mathbf{u}$  takes the value  $\mathbf{u} = (u_1, \dots, u_n)^\top$  with  $\phi$  being the (possibly vector-valued) parameter of the missingness process. It is possible for  $\theta$  to be a component of  $\phi$ , in which case the value of  $\theta$  itself plays a role in determining whether a response will be observed. In these cases we shall sometimes

write  $g(m | u, \theta, \phi)$  to emphasize the dependence on  $\theta$ .

•  $h_\theta(u | m)$  is the probability that  $\mathbf{u}$  takes the value  $u$  given that  $\mathbf{m}$  takes the value  $m$ .

•  $t_\phi(m)$  is the probability that  $\mathbf{m}$  takes the value  $m$ . Again,  $\theta$  may be a component of  $\phi$ .

### Example (continued)

Suppose that the missingness process in the two-item example initiated above can be described as follows: The second response is observed whenever the first response is correct; the second response will be observed with probability  $\phi$  if the first response is incorrect. Then

$$g_\phi(m | u) = \begin{cases} 1 & \text{if } m = (1, 0) \text{ and } u = (1, 0) \text{ or } (1, 1) \\ 1 - \phi & \text{if } m = (1, 0) \text{ and } u = (0, 0) \text{ or } (0, 1) \\ \phi & \text{if } m = (1, 1) \text{ and } u = (1, 0) \text{ or } (1, 1) \\ 0 & \text{otherwise.} \end{cases}$$

Whenever not all potential responses may be observed for any reason—even if they all do turn out to be observed—the data are  $\mathbf{v}$ . To obtain the likelihood function, we start with the likelihood for the (hypothetical) complete data  $\mathbf{u}, \mathbf{m}$ , then average over the missing responses  $u_{\text{mis}}$ :

$$L(\theta, \phi | \tilde{v}) = \delta[(\theta, \phi), \Omega_{\theta\phi}] \cdot \int f_\theta(u_{\text{mis}}, \tilde{u}_{\text{obs}}) g_\phi(\tilde{m} | u_{\text{mis}}, u_{\text{obs}}) du_{\text{mis}},$$

where  $\delta$  takes the value 1 if a value  $(\theta, \phi)$  is in the parameter space  $\Omega_{\theta\phi}$  and 0 if it is not. This observed-data likelihood is a weighted average over all complete-data likelihoods that have the targeted responses to the observed items. The weights are proportional to the probabilities of these potential response patterns for the different values  $u_{\text{mis}}$ , given  $\tilde{m}$  and  $\tilde{u}_{\text{obs}}$ . Using local independence, we can bring the probability for the observed responses outside the integral:

### References

- Little, R. J. A., & Rubin, D. B. (1987). *Statistical analysis with missing data*. Wiley.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3), 581–592. <https://doi.org/10.1093/biomet/63.3.581>