Measurement and Research Department Reports

On the Loss of Information in Conditional Maximum Likelihood Estimation of Item Parameters

T.J.H.M. Eggen



Measurement and Research Department Reports

97-6

On the Loss of Information in Conditional Maximum Likelihood Estimation of Item Parameters

T.J.H.M. Eggen

Cito Arnhem, 1997 Cite Instituut voor Toetsontwikkeling Postbus 1034 6801 MG Arnhem

Bibliotheek



This manuscript has been submitted for publication. No part of this manuscript may be copied or reproduced without permission.

Abstract

In item response models of the Rasch type (Fischer & Molenaar, 1995), item parameters are often estimated by the conditional maximum likelihood (CML) method. This paper addresses the loss of information in CML estimation by using the information concept of F-information (Liang, 1983). This concept makes it possible to specify the conditions for no loss of information and to define a quantification of information loss. For the dichotomous Rasch model, the derivations will be given in detail to show the use of the F-information concept for making efficiency comparisons for different estimation methods. It is shown that by using CML for item parameter estimation, some information is almost always lost. But compared to JML (joint maximum likelihood) as well as to MML (marginal maximum likelihood) the loss is very small. The reported efficiency of CML to JML and to MML in several comparisons is always larger than 93%, and in tests with a length of 20 items or more, larger than 99%.

Keywords: Conditional maximum likelihood, Loss of information, F-information, Rasch models, item parameter estimation

Introduction

The Rasch model for measuring a latent trait θ with dichotomously scored items, with the responses $X_{vi} = x_{vi}$ (0 or 1), persons v = 1, ..., n and items i = 1, ..., k is given by

$$p(x;\omega) = \prod_{\nu=1}^{n} \prod_{i=1}^{k} P(X_{\nu i} = x_{\nu i}; \beta_{i}, \theta_{\nu}), \qquad (1)$$

with
$$P(X_{\nu i} = x_{\nu i}; \beta_i, \theta_\nu) = \frac{\exp x_{\nu i}(\theta_\nu - \beta_i)}{1 + \exp(\theta_\nu - \beta_i)},$$
 (2)

and $\omega^T = (\beta^T, \theta^T)$, $\beta^T = (\beta_1, ..., \beta_k)$ is the vector item parameter (difficulty) and $\theta^T = (\theta_1, ..., \theta_n)$ the vector person parameter (ability).

In Fischer and Molenaar (1995) the foundations, the main results and an overview of recent developments and extensions of the Rasch model are given. In these type of item response theory (IRT) models conditional maximum likelihood estimation (CML) is a popular method for estimating item parameters. Although properties as consistency, asymptotic normality and sampling independence of CML estimators are well established, the justification of CML estimation with respect to the efficiency of the estimators is not clear. This topic is addressed in this paper.

Estimation of item parameters in the Rasch model

For the dichotomous Rasch model (1), three likelihood based methods for item parameter estimation are available (Molenaar, 1995).

In the first method, joint maximum likelihood (JML), the item parameters are estimated by maximizing (1) with respect to ω , given the data x. It is well known that estimating the item parameters by JML leads to inconsistent estimators (Andersen, 1973). This is caused by the fact that a limited number of interest (items) parameters (β) are to be estimated in the presence of many nuisance (ability) parameters (θ). Eliminating the nuisance parameters gives the solution for this problem (Basu, 1977). In IRT modeling this elimination is accomplished by the marginal or the conditional maximum likelihood method.

In the marginal maximum likelihood (MML) method it is assumed that the abilities θ_{y} in (1) constitute a random sample from an ability distribution $h(\theta; \xi)$, with ξ the

parameters of the ability distribution. In the Rasch model the joint probability of the item responses can then be written as

$$p_{m}(x;\beta,\xi) = \int p(x \mid \theta;\beta,\xi) h(\theta;\xi) d\theta =$$

$$\prod_{\nu=1}^{n} \int_{-\infty}^{\infty} \prod_{i=1}^{k} P(X_{\nu i} = x_{\nu i} \mid \theta_{\nu};\beta_{i}) h(\theta_{\nu};\xi) d\theta_{\nu} = L_{m}(\beta,\xi;x).$$
(3)

This so called marginal likelihood function is maximized with respect to β and ξ in order to get estimates of the item parameters, and, for instance, in the case of a normal ability distribution, two distribution parameters. In the MML method, the nuisance parameters are eliminated by integrating them out. In (3) $P(X_{vi} = x_{vi} | \theta_v; \beta_i)$ is given by (2).

In Rasch type of IRT models, the CML method is an alternative solution to the inconsistency problem. If in the model there exist sufficient statistics for the nuisance parameters, it can be separated in a conditional part which is only dependent on the interest parameters and a part which models the sufficient statistics.

In the dichotomous Rasch model the sum score on the items $T_{\nu} = \sum_{i=1}^{n} X_{\nu i}$ is sufficient for θ_{ν} , $\nu = 1, ..., n$, so we can rewrite (1) as

$$p(x;\omega) = \prod_{\nu=1}^{n} f(x_{\nu} \mid t_{\nu}; \beta) \cdot \prod_{\nu=1}^{n} g(t_{\nu}; \beta, \theta), \qquad (4)$$

with $X_{\nu} = (X_{\nu_1}, ..., X_{\nu_k})$ the response vector of person $\nu = 1, ..., n$. Maximizing, with respect to β , the conditional likelihood

$$L_{c}(\beta; x \mid t) = \prod_{\nu=1}^{n} f(x_{\nu} \mid t_{\nu}; \beta)$$
(5)

leads to consistent estimators of the item parameters (Andersen, 1973).

Information and efficiency

In likelihood inference, the information concept plays an important role. In the case of several parameters, the Fisher information matrix is used in the evaluation of the quality of estimators. It is well known (Rao, 1973) that in the class of unbiased estimators the variance of any estimator is bounded by the inverse of the Fisher information matrix. Estimators reaching this bound, the Cramér-Rao lower bound, are called efficient. And the efficiency of any other estimator is always expressed in relation to this lower bound.

It has been shown (Andersen, 1973) that the CML estimators for the item parameters in the Rasch model are under mild conditions asymptotically efficient; the same is true for the MML estimators of the item parameters. However, these asymptotic results do not imply that there is no loss of information if CML estimation is applied. The point is that by using the conditional likelihood (5) for estimating the item parameters, the second part of the full likelihood (4), which is the marginal distribution of T, is neglected. And this distribution possibly contains some information on the item parameters. In MML estimation, where no information is discarded, on the other hand, a correct specification of the ability distribution is needed and if this distribution is not the correct one, the resulting estimates of the item parameters can be useless. Furthermore, in MML the loss in efficiency of the item parameter estimation due to the joint estimation of them with the parameters of the ability distribution is not clear. In psychometric textbooks, the efficiency problem of the item parameter estimators is mostly discarded. In some recent work, however, Engelen (1989) and Zwinderman (1991), point out that there must be some loss of information on the item parameters without giving a quantification of the loss. Assuming that the ability distribution parameters are known, Engelen (1989) has shown that the loss in CML estimation compared to MML is small.

In this paper a general treatment is given of an information concept, called Finformation, which makes it possible to define a clear measure of information loss in using CML estimation. The conditions for no loss of information in separable models, like the Rasch model, will be given. The theory will be clarified with some examples. In particular, attention will be paid to the dichotomous Rasch model as specified in (1) and (2) and to the Rasch Poisson Counts model for misreadings (Rasch, 1960). For the dichotomous Rasch model the derivations will be given in detail to show the use of the general information concept for making efficiency comparisons for different estimation methods. Results will be given of the comparison of CML to JML, and of the comparison of CML to MML in case of a normal ability distribution.

Notation and Terminology

It should be noted that when the terms one- and two-parameter distribution are used in this paper, the parameters can be vectors. By a two-parameter distribution is meant that the parameters of the distribution ω can be partitioned in an interest parameter ψ and a nuisance parameter τ . In general, a scalar or vector random variable X is considered

with density (or probability in the discrete case) $p(x;\omega)$, with $\omega \in \Omega$. The parameter of interest $\psi = \psi(\omega)$ has domain Ψ and the nuisance parameter τ , the complementary part to ψ of ω , has domain T. It is assumed that $\Omega = \Psi \times T$. Let T = T(X) be a (vector) statistic, f(. | .;.) a conditional density or probability of the data given a statistic and g(.;.) the (marginal) density (or probability) of the statistic. Then in general the probability of the data can be factored in

$$p(x;\omega) = f(x \mid t;\omega) \cdot g(t;\omega)$$
(6)

In case T(X) is sufficient for τ then (6) can be written as

$$p(x;\omega) = f(x \mid t;\psi) \cdot g(t;\omega)$$
(7)

The model is then called separable and in CML estimation the inference is based on the first part of the factorization in (7) which is only dependent on the parameter of interest ψ , while the other part is neglected.

The F-Information: Definition and Basic Properties

F-information is a generalization of the Fisher information. In an *m*-dimensional oneparameter distribution $p(x:\omega)$, $\omega^T = (\omega_1, ..., \omega_m)^T$, the Fisher information matrix is defined as the $m \times m$ matrix

$$\mathbf{I}_{p}(\omega) = \mathscr{E}\left[S_{p;\omega}, S_{p;\omega}^{T}\right] , \qquad (8)$$

in which

$$S_{\rho;\omega} = S_{\rho;\omega}(X) := \frac{\partial \ln p(X;\omega)}{\partial \omega}, \qquad (9)$$

is the efficient score statistic (a $m \times 1$ -vector) of a distribution p with respect to ω .

The role of this Fisher information matrix in estimation problems will be illustrated by some examples. The Cramér-Rao lower bound for the covariance matrix of any estimator in the class of unbiased estimators is given by the inverse (if it exists) of the Fisher information matrix: if V is the covariance matrix of an unbiased estimator for ω , then $V - I_p^{-1}(\omega)$ is non-negative definite (nnd). Furthermore, it is known that if there exists an efficient estimator (reaching the lower bound), the maximum likelihood estimation (MLE) procedure will produce it. MLE's need not to be unbiased, of course, but in broad classes of models and under mild regularity conditions, these estimators are asymptotically normally distributed with covariance matrix the inverse of the Fisher information matrix. In practice, the so called *observed* Fisher information matrix is used, to contrast it with the *expected* Fisher information defined in (8) and (9). The observed Fisher information matrix is obtained by disregarding the expectation in (8) and evaluating it at the MLE $\omega = \hat{\omega}$. The inverse of the observed information matrix is used as an estimator of the covariance matrix of $\hat{\omega}$. Its diagonal elements are used for obtaining standard errors of the estimates and for constructing confidence intervals.

For two-parameter distributions, the concept of F-information, originating from Efron (1977), was defined by Liang (1983) for scalar parameters and was generalized to vector parameters by Bhapkar (1989).

Definition 1.

In a two-parameter distribution $p(x:\omega)$, with $\omega^T = (\psi^T, \tau^T)$, $\psi^T = (\psi_1, \dots, \psi_k)$ being the interest parameter and $\tau^T = (\tau_1, \dots, \tau_n)$ the nuisance parameter, the *F*-information for ψ in *p* is given by

$$I_{p}(\psi;\omega) = \underset{N}{\operatorname{Min}} \quad \mathscr{E}[m(N) \cdot m(N)^{T}], \text{ with } m(N) = S_{p;\psi} - N^{T} \cdot S_{p;\tau}.$$
(10)

In (10), N is a $n \times k$ -matrix with constants and Min is the minimal matrix M_1 in the class M of non-negative definite matrices which can be written as $\mathscr{E}[m(N) \cdot m(N)^T]$ for some N, that is for any $M \in \mathbb{M} : M - M_1$ is nnd.

Under general regularity conditions, which entail the justification of interchanging the order of differentiation and integration (taking expectations), the F-information has the following properties:

1. F-information is a generalization of the Fisher information. Partition the parameter vector in itself and an empty part $\omega^T = (\omega^T, \emptyset)$ and observe that

$$\mathbf{I}_{p}(\omega;\omega) = \underset{N}{\operatorname{Min}} \mathscr{E}\left[\left(S_{p;\omega} - N^{T}S_{p;\varnothing}\right) \cdot \left(S_{p;\omega} - N^{T}S_{p;\varnothing}\right)^{T}\right] = \mathscr{E}\left[S_{p;\omega} \cdot S_{p;\omega}^{T}\right] = \mathbf{I}_{p}(\omega)$$

2. F-information can be defined in terms of Fisher information.

It can be shown (Bhapkar, 1989) that if the Fisher information matrix (8) is positive definite and is rewritten as a partitioned matrix:

$$\mathbf{I}_{\rho}(\omega) = \mathscr{C}\left[S_{\rho;\omega}, S_{\rho;\omega}^{T}\right] = \mathscr{C}\left[\begin{array}{cc}S_{\rho;\psi}, S_{\rho;\psi}^{T}, S_{\rho;\psi}, S_{\rho;\tau}^{T}\\S_{\rho;\tau}, S_{\rho;\psi}^{T}, S_{\rho;\tau}, S_{\rho;\tau}^{T}\end{array}\right] = : \begin{bmatrix}\mathbf{I}_{11} & \mathbf{I}_{12}\\\mathbf{I}_{21} & \mathbf{I}_{22}\end{bmatrix}$$
(11)

and its inverse as

$$I_p^{-1}(\omega) = \begin{bmatrix} I^{11} & I^{12} \\ I^{21} & I^{22} \end{bmatrix},$$

then the F-information matrix for ψ in p is equal to

$$I_{p}(\psi;\omega) = I_{11} - I_{12}I_{22}^{-1}I_{21} = (I^{11})^{-1}, \qquad (12)$$

which is the inverse of the upper left submatrix of the inverse of the Fisher information matrix.

Suppose both ψ and τ are scalar, then $N = N^T$ and all efficient scores are scalar. The F-information is found by minimizing $\mathscr{C}(S_{p,\psi}^2 - 2NS_{p,\psi}S_{p,\tau} + N^2S_{p,\tau}^2)$ with respect to N. After differentiation, the minimum is easily seen to be reached at $N = \mathscr{C}(S_{p,\psi} S_{p,\tau}) \cdot (\mathscr{C}(S_{p,\tau})^2)^{-1}$, and the F-information is given by (12). Generalizing to the multidimensional case, the same result follows for: $N^T = \mathscr{C}(S_{p,\psi} S_{p,\tau}^T) \cdot \mathscr{C}(S_{p,\tau} S_{p,\tau}^T)^{-1}$.

3. F-information has a clear geometrical interpretation.

To enable a geometrical interpretation of the F-information observe: $I_p(\psi; \omega)$ is the minimal matrix in the class M of non-negative definite $k \times k$ matrices: $\forall K \in M : K - I_p(\psi; \omega) \in M$. Consider a real valued function U from $M \to \mathbb{R}$, so from the class of nnd matrices to the real line. Suppose this function satisfies:

a.
$$U(K) = 0$$
 if and only if $K = 0$ and (13a)

b. for every
$$K_1, K_2 \in M$$
 with $K_1 - K_2 \in M$ implies $U(K_1) \ge U(K_2)$. (13b)

Then it follows from (13b) that for every $K \in M$, by definition $I_{\rho}(\psi; \omega) \in M$ and $K - I_{\rho}(\psi; \omega) \in M$, which implies $U(K) \ge U(I_{\rho}(\psi; \omega))$.

So it can be concluded that finding the F-information is equivalent to finding the minimum of a function U. A matrix function which satisfies (13a) and (13b) is tr(.), the trace function. So minimizing $\mathscr{E}[S_{p;\psi} - N^T S_{p;\tau}) \cdot (S_{p;\psi} - N^T S_{p;\tau})^T]$, see (10), is equivalent to minimizing the trace of this matrix and then from standard matrix algebra it follows

$$\operatorname{Mintr}_{N} \left\{ \mathscr{E} \left[\left(S_{p;\psi} - N^{T} S_{p;\tau} \right) \cdot \left(S_{p;\psi} - N^{T} S_{p;\tau} \right)^{T} \right] \right\} = \\
\operatorname{Mintr}_{N} \left\{ \mathscr{E} \left[\left(S_{p;\psi} - N^{T} S_{p;\tau} \right)^{T} \cdot \left(S_{p;\psi} - N^{T} S_{p;\tau} \right) \right] \right\} = \\
\operatorname{Mintr}_{N} \mathscr{E} \left[\left(S_{p;\psi} - N^{T} S_{p;\tau} \right)^{T} \cdot \left(S_{p;\psi} - N^{T} S_{p;\tau} \right) \right] \right\}$$

The last expression makes it possible to give the F-information the same geometric interpretation as a least squares solution in a regression problem. As can be seen in Figure 1, the F-information follows from projecting the efficient score of the interest parameter onto the space of the efficient score of the nuisance parameter.



Figure 1 Geometrical interpretation F-information

4. F-information is the Cramér-Rao bound.

For a finite number of parameters, $I_p(\psi; \omega)$ is the Cramér-Rao lower bound in the class of unbiased estimators of ψ , that is, for every covariance matrix V of unbiased estimators of ψ V - $(I_p(\psi; \omega))^{-1}$ is non-negative definite (Rao, 1973, p. 326). So when an estimator of ψ in a two-parameter distribution $p(x; (\psi, \tau))$ is evaluated with respect to reaching this lower bound, it is seen that all the relevant information on ψ in this distribution is included in the F-information.

5. A scalar measure of information.

The F-information for a parameter ψ is in general a $k \times k$ -matrix. And comparing F-information matrices is not straightforward. However, following Bhapkar (1989, p. 147), a scalar measure of information can be defined by a function on the F-information matrix satisfying (13a) and (13b). Possible choices could be the largest eigenvalue of the F-information, the trace of the Finformation, and the determinant of the F-information in case the matrix is nonsingular. In this paper only the trace of the F-information matrix, $tr(I_p(\psi; \omega))$, will be used as a scalar measure of information on ψ in p.

6. F-information in case Fisher information is not positive definite.

Bhapkar (1991) remarks that even in case the Fisher information matrix is not positive definite the F-information exists, although (12) cannot be used directly. In general, it yields that for any $n \times k$ matrix B which is a solution of $I_{22} B = I_{21}$ the F-information is given by $I_p(\psi; \omega) = I_{11} - B^T I_{22} B$, which is the same for any B.

Note that in the special case that the submatrix I_{22} of $I_p(\omega)$ is positive definite, a solution is given by $B = I_{22}^{-1} I_{21}$, and it follows that

$$\mathbf{I}_{p}(\psi;\omega) = \mathbf{I}_{11} - \mathbf{I}_{21}^{T} \left(\mathbf{I}_{22}^{-1}\right)^{T} \mathbf{I}_{22} \ \mathbf{I}_{22}^{-1} \ \mathbf{I}_{21} = \mathbf{I}_{11} - \mathbf{I}_{12} \ \mathbf{I}_{22}^{-1} \ \mathbf{I}_{21}, \tag{14}$$

which gives the same expression as (12).

F-Information in Separable Models

The F-information can be used to quantify the loss of information in estimating the interest parameters ψ in case the original two-parameter model can be factored as in (7), $p(x;\omega) = f(x \mid t; \psi) \cdot g(t;\omega)$, by just considering the conditional distribution $f(x \mid t; \psi)$. Moreover, conditions can be specified for no loss of information using this method.

In these separable models, some properties of the F-information will be shown to hold. Throughout it is assumed, that regularity conditions are met which guarantee the existence of the Fisher information and allow interchanging the order of differentiation and integration of the logarithm of the model.

Properties of the efficient score statistics

Three basic properties of the efficient score statistics, involved in the definition of Finformation, are given first.

Property 1

$$\mathscr{E} S_{p;\psi} = \mathscr{E} S_{p;\tau} = 0.$$

Property 2
a.
$$\operatorname{COV}(S_{p;\psi}, S_{p;\psi}^T) = \mathscr{C}(S_{p;\psi}S_{p;\psi}^T) = -\mathscr{C}S_{p;\psi,\psi^T} := -\mathscr{C}\left(\frac{\partial^2 \log p}{\partial \psi \partial \psi^T}\right);$$
 and
b. $\operatorname{COV}(S_{p;\tau}, S_{p;\tau}^T) = \mathscr{C}(S_{p;\tau}S_{p;\tau}^T) = -\mathscr{C}S_{p;\tau,\tau^T}.$

Property 3
a.
$$\operatorname{COV}(S_{p;\psi}, S_{p;\tau}^T) = \mathscr{C}(S_{p;\psi}S_{p;\tau}^T) = -\mathscr{C}S_{p;\psi,\tau^T} := -\mathscr{C}\left(\frac{\partial^2 \log p}{\partial \psi \, \partial \tau^T}\right)$$
; and
b. $\operatorname{COV}(S_{p;\tau}, S_{p;\psi}^T) = \mathscr{C}(S_{p;\tau}S_{p;\psi}^T) = -\mathscr{C}S_{p;\tau,\psi^T}.$

These three properties are valid in each two-parameter distribution. So when a decomposition as in (6) or (7) is considered, the properties also apply to the efficient score statistics of $f(x \mid t; \omega)$ and $g(t; \omega)$. Then the following properties, relating the efficient score statistics of p, f and g, are easily deduced.

Property 4 a. $\mathscr{E}(S_{p;\psi} \mid T = t) = S_{g;\psi}$; and b. $\mathscr{E}(S_{p;\tau} \mid T = t) = S_{g;\tau}$.

Proof 4a

$$\mathscr{E}(S_{p;\psi} \mid T=t) = \int_{\{x:t(x)=t\}} \frac{\partial \log p(x;\omega)}{\partial \psi} f(x \mid t;\omega) dx = \int_{\{x:t(x)=t\}} \frac{\partial p(x;\omega)/\partial \psi}{p(x;\omega)} \frac{\partial p(x;\omega)}{g(t;\omega)} dx = \int_{\{x:t(x)=t\}} \frac{\partial p(x;\omega)/\partial \psi}{g(t;\omega)} dx = \frac{\partial g(t;\omega)/\partial \psi}{g(t;\omega)} = \frac{\partial \log g(t;\omega)}{\partial \psi} = S_{g;\psi}.$$

The fourth equality yields because $g(t;\omega) = \int_{\{x:t(x)=t\}} p(x;\omega) dx$. \Box

Property 5 a. $\mathscr{E}\left(S_{p;\psi}S_{g;\psi}^{T}\right) = \mathscr{E}\left(S_{g;\psi}S_{g;\psi}^{T}\right);$ and b. $\mathscr{E}\left(S_{p;\tau}S_{g;\tau}^{T}\right) = \mathscr{E}\left(S_{g;\tau}S_{g;\tau}^{T}\right).$

Proof 5a

Following a general property of conditional expectations it yields

$$\mathscr{E}\left(S_{p;\psi}S_{g;\psi}^{T}\right) = \mathscr{E}_{t}\left(\mathscr{E}_{x\mid t}\left(S_{p;\psi}S_{g;\psi}^{T}\mid T=t\right)\right),$$

but in the conditional distribution $S_{g;\psi}$ is a constant. So

$$\mathscr{E}_{x|t}\left(S_{p;\psi} \; S_{g;\psi}^T \mid T=t\right) = \mathscr{E}_{x|t}\left(S_{p;\psi} \mid T=t\right) \; S_{g;\psi}^T.$$

Then using property 4a, 5a follows. \Box

In the special case of a separable two-parameter distribution that can be decomposed as in (7), the following properties can be added.

Property 6

a. $S_{\rho;\psi} = S_{f;\psi} + S_{g;\psi}$, by definition.

The same yields for score statistics with respect to τ . But in this case:

b.
$$S_{f,\tau} = 0$$
 and $S_{p;\tau} = S_{g;\tau}$.

Property 7
a.
$$\mathscr{E}\left(S_{f;\psi}S_{g;\tau}^{T}\right) = 0$$
; and

b.
$$\mathscr{E}\left(S_{f;\tau}S_{g;\psi}^{T}\right)=0.$$

Proof 7a

Because of property 3 we have:

$$\mathscr{E}\left(S_{p;\psi}S_{p;\tau}^{T}+S_{p;\psi,\tau}\right)=\mathscr{E}\left(S_{p;\psi}S_{p;\tau}^{T}+\frac{\partial}{\partial\psi}S_{p;\tau}^{T}\right)=0.$$

Applying property 6 this can be rewritten as:

$$\begin{split} & \mathscr{E}\left[\left(S_{f;\psi}+S_{g;\psi}\right)S_{g;\tau}^{T}+\frac{\partial}{\partial\psi}S_{p;\tau}^{T}\right]=0, \\ & \mathscr{E}\left[S_{f;\psi}S_{g;\tau}^{T}+S_{g;\psi}S_{g;\tau}^{T}+\frac{\partial}{\partial\psi}S_{g;\tau}^{T}\right]=0, \\ & \mathscr{E}\left(S_{f;\psi}S_{g;\tau}^{T}\right)+\mathscr{E}\left(S_{g;\psi}S_{g;\tau}^{T}\right)+\mathscr{E}S_{g;\psi,\tau}^{T}=0. \end{split}$$

And the last two terms in the left hand side of the last equation sum to 0, because of property $3a. \square$

Property 8 a. $\mathscr{C}\left(S_{f;\psi}S_{g;\psi}^{T}\right) = 0$, and b. $\mathscr{C}\left(S_{f;\tau}S_{g;\tau}^{T}\right) = 0$.

Proof 8a Because of property 5a

$$\mathscr{E}\left(S_{p;\psi} S_{g;\psi}^{T} - S_{g;\psi} S_{g;\psi}^{T}\right) = 0,$$
$$\mathscr{E}\left[\left(S_{p;\psi} - S_{g;\psi}\right) S_{g;\psi}^{T}\right] = 0.$$

Using property 6 this becomes

$$\mathscr{E}\left[\left(S_{f;\psi} + S_{g;\psi} - S_{g;\psi}\right) S_{g;\psi}^{T}\right] = \mathscr{E}\left(S_{f;\psi} S_{g;\psi}^{T}\right) = 0. \Box$$

Theorems on the F-information in separable models

Theorem 1

For a two-parameter distribution which can be factored in $p(x;\omega) = f(x \mid t; \psi) \cdot g(t;\omega)$

a. the F-information for ψ in p is the sum of the F-information in f and the F-information in g:

$$I_{p}(\psi;\omega) = I_{f}(\psi;\omega) + I_{p}(\psi;\omega).$$
(15a)

b. Moreover, the F-information in the conditional distribution is the same as the expectation of the Fisher information in the conditional distribution, with the expectation taken with respect to the distribution of the sufficient statistic T:

$$I_{f}(\psi;\omega) = \mathscr{E}I_{f}(\psi \mid T). \text{ So, } I_{p}(\psi;\omega) = \mathscr{E}I_{f}(\psi \mid T) + I_{g}(\psi;\omega).$$
(15b)

Proof Theorem 1:

a. Using (11) and (12) we write

$$\mathbf{I}_{p}(\psi;\omega) = \mathscr{E}\left[S_{p;\psi}, S_{p;\psi}^{T}\right] - \mathscr{E}\left[S_{p;\psi}, S_{p;\tau}^{T}\right] \left(\mathscr{E}\left[S_{p;\tau}, S_{p;\tau}^{T}\right]\right)^{-1} \mathscr{E}\left[S_{p;\tau}, S_{p;\psi}^{T}\right].$$
(16)

Replacing in (16) the efficient score statistics in p by those in f and g, using property 6 above, gives:

$$\begin{split} \mathbf{I}_{p}(\psi;\omega) &= \mathscr{C}\left(S_{f;\psi}S_{f;\psi}^{T} + S_{f;\psi}S_{g;\psi}^{T} + S_{g;\psi}S_{f;\psi}^{T} + S_{g;\psi}S_{g;\psi}^{T}\right) - \\ \mathscr{C}\left(S_{f;\psi}S_{g;\tau}^{T} + S_{g;\psi}S_{g;\tau}^{T}\right) \left(\mathscr{C}\left[S_{g;\tau}S_{g;\tau}^{T}\right]\right)^{-1} \mathscr{C}\left(S_{g;\tau}S_{f;\psi}^{T} + S_{g;\tau}S_{g;\psi}^{T}\right) \end{split}$$

In this expression, due to properties 7 and 8,

$$\mathscr{E}\left(S_{f;\psi}S_{g;\psi}^{T}\right) = \mathscr{E}\left(S_{g;\psi}S_{f;\psi}^{T}\right) = \mathscr{E}\left(S_{f;\psi}S_{g;\tau}^{T}\right) = \mathscr{E}\left(S_{g;\tau}S_{f;\psi}^{T}\right) = 0. \text{ So,}$$
$$\mathbf{I}_{p}(\psi;\omega) = \mathscr{E}\left(S_{f;\psi}S_{f;\psi}^{T}\right) + \mathscr{E}\left(S_{g;\psi}S_{g;\psi}^{T}\right) - \mathscr{E}\left(S_{g;\psi}S_{g;\tau}^{T}\right) \left(\mathscr{E}\left[S_{g;\tau}S_{g;\tau}^{T}\right]\right)^{-1} \mathscr{E}\left(S_{g;\tau}S_{g;\psi}^{T}\right).$$

In the right hand side of this expression, according to (12), the last two terms give the F-information in g, and $\mathscr{C}(S_{f;\psi}S_{f;\psi}^T) = I_f(\psi;\omega)$, the F-information in f, because $S_{f,\tau} = 0$. This results in:

$$\mathbf{I}_{p}(\psi;\omega) = \mathbf{I}_{f}(\psi;\omega) + \mathbf{I}_{g}(\psi;\omega).\Box$$

b. From the definitions and a property of conditional expectations follows:

$$I_{f}(\psi;\omega) = \mathscr{C}_{x}\left[S_{f;\psi}, S_{f;\psi}^{T}\right] = \mathscr{C}_{t}\left[\mathscr{C}_{x\mid t}\left(S_{f;\psi}, S_{f;\psi}\mid T\right)\right] = \mathscr{C}I_{f}(\psi\mid T), \text{ which gives:}$$

$$I_{p}(\psi;\omega) = \mathscr{C}I_{f}(\psi\mid T) + I_{g}(\psi;\omega).\Box$$

In separable two-parameter distributions therefor, the F-information with respect to the interest parameter ψ can be written as the sum of the F-information in the conditional distribution and the F-information in the marginal distribution of the statistic T(X) This additivity of the components of the F-information is a very attractive. It will be clear that the loss of information by using the conditional model instead of the full two-parameter model in an inference on ψ is given by the F-information in T: $I_g(\psi; \omega)$. The efficiency of using the conditional instead of the full model can be expressed as $I_f(\psi; \omega)/I_p(\psi; \omega)$. However, it should be remembered, that in general, having a k-vector interest parameter ψ , the F-informations are $k \times k$ matrices. For comparing matrices and computing efficiencies, the trace function of the matrices will be used. The efficiency of two models with respect to the information on ψ will be computed as:

$$EFF(\psi;f:p) = \frac{\operatorname{tr}\left(I_{f}(\psi;\omega)\right)}{\operatorname{tr}\left(I_{p}(\psi;\omega)\right)}$$
(17)

Example

An example of the computation of the F-information is given for a problem in which both the interest and the nuisance parameter are scalar. Consider a random sample from a normal distribution $X = (X_1, ..., X_n)$, $X_i \sim N(\mu, \sigma^2)$, with both the mean μ and the variance σ^2 unknown. The interest is in estimating σ^2 , while μ is considered as a nuisance parameter. In this problem, the sample mean $T = \sum_i X_i / n$ is sufficient for μ and the full two-parameter model can be decomposed as in (7):

$$p(x;\mu,\sigma^2) = f(x \mid t;\sigma^2) g(t;\mu,\sigma^2).$$

For estimating σ^2 , the full model p or the conditional model f can be used.

The Fisher information matrices for the full model p and the distribution g of T, $(T \sim N(\mu, \sigma^2/n))$, are

$$I_{p}(\mu, \sigma^{2}) = \begin{bmatrix} n\sigma^{-2} & 0 \\ 0 & \frac{n\sigma^{-4}}{2} \end{bmatrix}, \text{ and } I_{g}(\mu, \sigma^{2}) = \begin{bmatrix} n\sigma^{-2} & 0 \\ 0 & \frac{\sigma^{-4}}{2} \end{bmatrix}.$$

The F-information with respect to σ^2 in p and g follow from (12):

 $I_p(\sigma^2; (\mu, \sigma^2)) = n/2\sigma^4$, and $I_g(\sigma^2; (\mu, \sigma^2)) = 1/2\sigma^4$. From Theorem 1 (15a), the Finformation in the conditional model equals: $I_f(\sigma^2; (\mu, \sigma^2)) = (n-1)/2\sigma^4$. So in using CML instead of ML for estimating σ^2 , there is a small loss of information which amounts to $1/2\sigma^4$. The efficiency of CML versus ML is given by: $I_f(\sigma^2; (\mu, \sigma^2)) / I_p(\sigma^2; (\mu, \sigma^2)) = (n-1)/n$. In this example it can be remarked that, although there is a little loss of information in using CML, the resulting CML estimator is unbiased, whereas the ML estimator is not: $\hat{\sigma}_{cml}^2 = \sum_i (x_i - t)^2/(n-1)$ and $\hat{\sigma}_{ml}^2 = \sum_i (x_i - t)^2/n$. \Box

The conditions for no loss of information

Next, two definitions are given which make it possible to specify the conditions under which there is no loss of information when the conditional model is used instead of the full model.

Definition 2.

In a two-parameter distribution, satisfying (7), the statistic T = T(X) is strongly ancillary for ψ , if $g(t;\omega)$ only depends on $\tau: p(x;\omega) = f(x \mid t;\psi) g(t;\tau)$.

It is easily understood that in case the statistic T is strongly ancillary for ψ , there will be no loss of information if in an inference on ψ the conditional distribution $f(x \mid t; \psi)$ is used instead of $p(x; \omega)$. In this case the part which is neglected does not depend on the parameter ψ and $I_g(\psi; \omega) = 0$.

Bhapkar (1989), generalizing an earlier result of Liang (1983) for scalar parameters to vector parameters, has proven that there is also no loss of information if weaker conditions on the distribution of T are fulfilled. Consider the following definition:

Definition 3.

In a two-parameter distribution, satisfying (7), the statistic T = T(X) is weakly ancillary for ψ , if there exists a one-to-one reparametrization between (ψ, τ) and (ψ, δ) such that $g(t;\omega) = g(t;\delta)$, that is, only depends on δ .

Example

 X_1 and X_2 are independent Poisson distributed variables: $X = (X_1, X_2), \mathscr{E}(X_1) = \lambda$ and $\mathscr{E}(X_2) = \lambda \rho$. The interest and nuisance parameter are respectively $\ln \rho$ and $\ln \lambda$. $T(X) = X_1 + X_2$ is also Poisson distributed with $\mathscr{E}(T) = \lambda + \lambda \rho$. Factorization (7) yields:

$$p(x;\lambda,\rho) = \frac{\exp\left[\left(x_1 + x_2\right)\ln\lambda + x_2\ln\rho - \left(\lambda + \lambda\rho\right)\right]}{x_1! x_2!} =$$

$$\frac{t! \exp[x_2 \ln\rho]}{x_1! x_2! \exp[t \ln(1+\rho)]} \cdot \frac{\exp[t \ln\lambda + t \ln(1+\rho) - (\lambda+\lambda\rho)]}{t!} = f(x \mid t;\rho) \cdot g(t;\lambda,\rho).$$

Apply the one-one reparametrization $\ln \rho = \ln \rho$ and $\ln \lambda = \ln \delta - \ln(1 + \rho)$. The distribution of T is then given by: $g(t;\rho,\delta) = (\exp[t\ln\delta - \delta])/t!$, which only depends on δ .

Next, Bhapkar's theorem (1989) is given.

Theorem 2

A two-parameter distribution which can be factored in $p(x;\omega) = f(x \mid t;\psi) g(t;\omega)$ and in which the statistic T = T(X) is weakly ancillary for ψ , there is no loss of information in using the conditional distribution $f(x \mid t;\psi)$ instead of $p(x;\omega)$ for inference on ψ : $I_{g}(\psi;\omega) = 0$.

The condition of weak ancillarity, although slightly differently defined, was also the key condition under which Andersen (1970, 1973) established the asymptotic efficiency of conditional maximum likelihood estimators. According to Andersen (1973, p. 99), weak ancillarity of the statistic T means that no inference about the interest parameter ψ can be drawn from the distribution of T, $g(t;\omega)$, that is not completely dependent on the specification of the nuisance parameter τ . In other words, nothing can be learned from the data about the interest parameter from the sole observation of the statistic. Although intuitively weak ancillarity may be an appealing concept, it is not easy to show in general that a distribution has this property. However, for models belonging to the exponential family, necessary and sufficient conditions for weak ancillarity, which are easily checked, are given in the next theorem (Andersen, 1973; Liang, 1983; Bhapkar 1989).

Theorem 3

If X has a distribution belonging to the exponential family with natural parameters $\omega = (\psi, \tau)$:

$$p(x;\omega) = c(\psi, \tau) \cdot \exp\left\{ u^{T}(x) \cdot \psi + t^{T}(x) \cdot \tau + b(x) \right\}$$
(18)

and the distribution of T = T(X) is given by

$$g(t;\omega) = c(\psi,\tau) \cdot \gamma(t,\psi) \cdot \exp(t^T \tau), \qquad (19)$$

then T is weakly ancillary for ψ if and only if there exist functions of ψ only and independent of the data: $w_i(\psi)$, i = 1, ..., k, and $v(\psi)$, such that:

$$\frac{\partial \ln \gamma(t,\psi)}{\partial \psi_i} = w_i(\psi) t + v(\psi), \text{ for } i = 1, \dots, k$$
(20)

for all (ae) t.

Checking the Conditions for no Loss of Information using CML in Two Rasch Models

The Rasch Poisson Counts model

In this model, proposed by Rasch (1960), the number of failures X_{vi} of person v = 1, ..., n on test i = 1, ..., k, which each consist of a number of items with low error probabilities, is considered. A well known application of the model is the number of misreadings in texts. The model assumes that X_{vi} is Poisson distributed with parameter $\lambda_{vi} = \beta_i \theta_v$, β_i being the difficulty parameter of text *i* and θ_v the ability parameter of person *v*. Writing the model in the exponential family form (18) gives:

$$\mathbf{P}(X_{\nu i}=x_{\nu i};\beta_i\theta_{\nu})=\frac{\exp(-\beta_i\theta_{\nu})\cdot\exp\{x_{\nu i}\ln\beta_i+x_{\nu i}\ln\theta_{\nu}\}}{x_{\nu i}!},\ x_{\nu i}=0,1,2,\ldots.$$

With the assumption of independence over texts and persons we get:

$$p(x;\beta,\theta) = \prod_{i} \prod_{\nu} P(X_{\nu i} = x\nu i;\beta_{i},\theta_{\nu}) =$$

$$\frac{\exp\{\sum_{i} \sum_{\nu} -\beta_{i}\theta_{\nu}\} \cdot \exp\{\sum_{i} \left(\sum_{\nu} x_{\nu i}\right)\ln\beta_{i} + \sum_{\nu} \left(\sum_{i} x_{\nu i}\right)\ln\theta_{\nu}\}}{\prod_{i} \prod_{\nu} x_{\nu i}!}.$$

It is easily seen that the model is a member of the exponential family (18), in which the statistic $T_v = \sum_i X_{vi}$, the number of failures of person v is sufficient for $\ln \theta_v$, for v = 1, ..., n.

So, the model is separable:

$$p(x;\beta,\theta) = \prod_{\nu=1}^{n} f(x_{\nu} \mid t_{\nu};\beta) \cdot \prod_{\nu=1}^{n} g(t_{\nu};\beta,\theta_{\nu}),$$

and for estimating the text parameters β CML estimation can be considered to use. Instead of the full model $p(x; \beta, \theta)$, only the conditional model, the first part of the factorization, is used. Whether neglecting the distribution of T causes loss of information is now easily checked.

Observe that T_{ν} is also Poisson distributed, with parameter $\theta_{\nu} \Sigma_{i} \beta_{i}$:

$$g(t;\beta,\theta) = \prod_{\nu=1}^{n} g(t_{\nu};\beta,\theta_{\nu}) =$$

$$\frac{\exp\left\{\sum_{i}\sum_{\nu} -\beta_{i}\theta_{\nu}\right\} \cdot \exp\left\{\left(\sum_{\nu} t_{\nu}\right)\ln\sum_{i}\beta_{i} + \sum_{\nu} t_{\nu}\ln\theta_{\nu}\right\}}{\prod_{\nu} t_{\nu}!}.$$

So the function $\gamma(t, \psi)$ in (19) is given by

$$\gamma(t;\beta) = \prod_{\nu=1}^{n} \left[\exp\left\{ t_{\nu} \ln \sum_{i} \beta_{i} \right\} \right] / t_{\nu}!,$$

and because $\ln \gamma(t;\beta) = \sum_{\nu} \left\{ t_{\nu} \ln \sum_{i} \beta_{i} - \ln (t_{\nu}!) \right\}$ $\frac{\partial \ln \gamma(t;\beta)}{\partial \beta_{j}} = \sum_{\nu} t_{\nu} \frac{1}{\sum_{i} \beta_{i}}, \text{ for } j = 1, \dots, k.$

The condition, (20), for T being weakly ancillary for β is seen to be fulfilled: $\partial \ln \gamma(t;\beta) / \partial \beta_j$ is a linear function of t, and $w_j(\psi) = 1 / \sum_i \beta_i$, j = 1, ..., k only depends on the text parameter β . So, there will be no loss of information if the text parameters are estimated with CML. In terms of F-information, with $\omega^T = (\beta^T, \theta^T)$: $I_p(\beta; \omega) = I_f(\beta; \omega)$ or $I_g(\beta; \omega) = 0$.

The Rasch Model for Dichotomously Scored Items

The same conditions will be checked for the model which was presented in the introduction. Writing this model, (1) and (2), as

$$p(x;\beta,\theta) = \frac{\exp\left\{-\sum_{i}\left(\sum_{\nu} x_{\nu i}\right)\beta_{i} + \sum_{\nu}\left(\sum_{i} x_{\nu i}\right)\theta_{\nu}\right\}}{\prod_{i}\prod_{\nu}\left\{1 + \exp\left(\theta_{\nu} - \beta_{i}\right)\right\}},$$
(21)

it can be seen to belong to the exponential family. $T_{\nu} = \sum_{i} X_{\nu i}$ is sufficient for θ_{ν} , for $\nu = 1, ..., n$. The distribution of T is checked for weak ancillarity. This distribution is given by

$$g(t;\beta,\theta) = \prod_{\nu=1}^{n} g(t_{\nu};\beta,\theta_{\nu}) = \prod_{\nu=1}^{n} \frac{\exp\left(\theta_{\nu}t_{\nu}\right) \cdot \gamma_{t_{\nu}}(\beta)}{\prod_{i=1}^{k} \left\{1 + \exp\left(\theta_{\nu} - \beta_{i}\right)\right\}},$$
(22)

with
$$\gamma_{t_{\nu}}(\beta) = \sum_{\sum x_{\nu}=t_{\nu}} \exp\left[-\sum_{i=1}^{k} \beta_{i} x_{\nu i}\right]$$

For, j = 1, ..., k, this gives

$$\frac{\partial \ln \prod_{\nu=1}^{n} \gamma_{t_{\nu}}(\beta)}{\partial \beta_{j}} = \sum_{\nu=1}^{n} \frac{\partial \ln \gamma_{t_{\nu}}(\beta)}{\partial \beta_{j}} = \sum_{\nu=1}^{n} \frac{e^{-\beta_{j}} \cdot \gamma_{t_{\nu}-1}^{(j)}}{\gamma_{t_{\nu}}(\beta)} \cdot \frac{t_{\nu}}{t_{\nu}} = \sum_{\nu=1}^{n} \frac{e^{-\beta_{j}} \cdot \gamma_{t_{\nu}-1}^{(j)}}{\sum_{i=1}^{k} e^{-\beta_{i}} \cdot \gamma_{t_{\nu}-1}^{(i)}} \cdot t_{\nu}, \quad (23)$$

in which

$$\gamma_{t_{v}-1}^{(j)} := \partial \gamma_{t_{v}}(\beta) / \partial e^{-\beta_{j}}, \text{ for } j = 1, \dots, k.$$

Note that the last equality in (23) uses an expression given by Fischer (1974, p. 242):

$$\gamma_{t_{v}}(\beta) . t_{v} = \sum_{i=1}^{k} e^{-\beta_{i}} \gamma_{t_{v}-1}^{i}.$$

It is seen that condition (20) of Theorem 3 is not fulfilled, since the function $w_j(\psi)$ is not only dependent on the item parameters β , but also on the statistic t_{ν} . T is therefore not weakly ancillary for β , which means that using CML estimation is possibly accompanied with loss of information compared to the situation in which the full model is used. In terms of F-information: $I_p(\beta; \omega) \ge I_f(\beta; \omega)$ or $I_g(\beta; \omega) \ge 0$.

F-Information in the Dichotomous Rasch Model: Comparing JML and CML

In order to determine the loss of information in the Rasch model, with both the itemand ability parameters considered fixed, the expressions for the F-information in the full model p, the distribution g of the sufficient statistic and conditional distribution f, are given first. The Fisher information matrix of the model $p(x;\beta,\theta)$, given in (21), with respect to both parameters β and θ , is written in the partitioned form as in (11):

$$\mathbf{I}_{p}(\omega) = \mathscr{E} \begin{bmatrix} S_{p;\beta} \ S_{p;\beta}^{T} \ S_{p;\beta} \ S_{p;\theta}^{T} \\ S_{p;\theta} \ S_{p;\beta}^{T} \ S_{p;\theta} \ S_{p;\theta}^{T} \end{bmatrix} =: \begin{bmatrix} \mathbf{I}_{p}^{\beta\beta^{T}} \ \mathbf{I}_{p}^{\beta\theta^{T}} \\ \mathbf{I}_{p}^{\theta\beta^{T}} \ \mathbf{I}_{p}^{\theta\theta^{T}} \end{bmatrix}$$
(24)

Using the properties 1, 2 and 3 of the efficient score statistics, the submatrices, $I_p^{\beta\beta^r}(k \times k)$, $I_p^{\beta\theta^r} = [I_p^{\theta\beta^r}]^T(k \times n)$ and $I_p^{\theta\theta^r}(n \times n)$, which determine the F-information (see (12)), are easily seen to be the negative expectation of the second derivatives of $\ln p$ with respect to the parameters. The elements of the submatrices in (24) are then, writing

$$p_{\nu i} := \frac{\exp(\theta_{\nu} - \beta_i)}{1 + \exp(\theta_{\nu} - \beta_i)},$$
(25)

given by:

$$\left(\mathbf{I}_{p}^{\beta\beta^{T}}\right)_{ii} = \sum_{\nu=1}^{n} p_{\nu i} \left(1 - p_{\nu i}\right) \text{ for } i = 1, \dots, k;$$
(26a)

$$\left(\mathbf{I}_{p}^{\beta\beta^{T}}\right)_{ij} = 0 \text{ for } i \neq j = 1, \dots, k;$$
(26b)

$$\left(I_{\rho}^{\beta\theta^{\tau}}\right)_{i\nu} = -p_{\nu i}(1-p_{\nu i}) \text{ for } i=1,...,k; \nu=1,...,n;$$
(26c)

$$\left(I_{\rho}^{\theta\theta^{r}}\right)_{\nu\nu} = \sum_{i=1}^{k} p_{\nu i} \left(1 - p_{\nu i}\right) \text{ for } \nu = 1, \dots, n;$$
(26d)

$$\left(I_{\rho}^{\theta\theta^{r}}\right)_{vw} = 0 \text{ for } v \neq w = 1, \dots, n.$$
(26e)

The Fisher information matrix is clearly singular. For every row and column

$$\sum_{i=1}^{k+n} (\mathbf{I}_p(\omega))_{ij} = \sum_{j=1}^{k+n} (\mathbf{I}_p(\omega))_{ij} = 0 \text{ for } i, j = 1, \dots, k + n.$$

But because the diagonal submatrix $I_p^{\theta\theta^r}$ is positive definite, and the elements of the inverse are given by the reciprocals of (26d), using (14), the F-information in the full Rasch model $I_p(\beta;\omega)$, is given by

$$\left(\mathbf{I}_{p}(\beta;\omega)\right)_{ii} = \sum_{\nu=1}^{n} p_{\nu i} \left(1 - p_{\nu i}\right) - \sum_{\nu=1}^{n} \frac{\left[p_{\nu i} \left(1 - p_{\nu i}\right)\right]^{2}}{\sum_{\ell=1}^{k} p_{\nu \ell} \left(1 - p_{\nu \ell}\right)} \text{ for } j = 1, \dots, k,$$
(27a)

$$\left(\mathbf{I}_{p}(\beta;\omega)\right)_{ij} = -\sum_{\nu=1}^{n} \frac{p_{\nu i}(1-p_{\nu i}) \cdot p_{\nu j}(1-p_{\nu j})}{\sum_{i=1}^{k} p_{\nu \ell}(1-p_{\nu \ell})} \text{ for } i \neq j = 1, \dots, k.$$
(27b)

Next, the F-information in the distribution $g(t;\beta,\theta)$ of the sufficient statistic (see (22)) is determined. The Fisher information matrix, analogous to (24), is given by

$$\mathbf{I}_{g}(\omega) = \mathscr{C}\begin{bmatrix} S_{g;\beta} S_{g;\beta}^{T} & S_{g;\beta} S_{g;\theta}^{T} \\ S_{g;\theta} S_{g;\beta}^{T} & S_{g;\theta} S_{g;\theta}^{T} \end{bmatrix} =: \begin{bmatrix} \mathbf{I}_{g}^{\beta\beta^{T}} & \mathbf{I}_{g}^{\beta\theta^{T}} \\ \mathbf{I}_{g}^{\theta\beta^{T}} & \mathbf{I}_{g}^{\theta\theta^{T}} \end{bmatrix}.$$
(28)

Because of the properties 3 and 6 $(S_{g,\theta} = S_{p,\theta})$ of the efficient score statistics, $I_g^{\beta\theta^r} = I_p^{\beta\theta^r}$ and $I_g^{\theta\theta^r} = I_p^{\theta\theta^r}$, which are given by (26c), (26d) and (26e). This is, of course, a consequence of the fact that T_v is sufficient for θ_v in $p(x; \beta, \theta)$. The second derivatives of $\ln g$ with respect to β are given by:

$$\partial^{2} \ln g / \partial \beta_{i}^{2} = \sum_{\nu=1}^{n} \left\{ \frac{e^{-\beta_{i}} \cdot \gamma_{t_{\nu}^{-1}}^{(i)}}{\gamma_{t_{\nu}}} - \left[\frac{e^{-\beta_{i}} \cdot \gamma_{t_{\nu}^{-1}}^{(i)}}{\gamma_{t_{\nu}}} \right]^{2} - p_{\nu i} (1 - p_{\nu i}) \right\} \text{ for } i = 1, \dots, k,$$

and

$$\frac{\partial^2 \ln g}{\partial \beta_j \partial \beta_i} = \sum_{\nu=1}^n \left\{ \frac{e^{-\beta_i - \beta_j} \cdot \gamma_{t_\nu - 2}^{(i,j)}}{\gamma_{t_\nu}} - \frac{e^{-\beta_i - \beta_j} \cdot \gamma_{t_\nu - 1}^{(i)} \gamma_{t_\nu - 1}^{(j)}}{\gamma_{t_\nu}^2} \right\} \text{ for } i \neq j = 1, \dots, k,$$

with

$$\gamma_{t_{v}-2}^{(i,j)} := \frac{\partial^{2} \gamma_{t_{v}}(\beta)}{\partial e^{-\beta_{i}} \partial e^{-\beta_{j}}}.$$

Observe that in the dichotomous Rasch model

$$P(X_{vi} = 1 \mid T_v = t_v; \beta) = \frac{e^{-\beta_i} \cdot \gamma_{t_v-1}^{(i)}}{\gamma_{t_v}} =: p_{vi \mid t_v} \text{ for } i = 1, ..., k$$

$$\mathbf{P}(X_{\nu i}=1, X_{\nu j}=1 \mid T_{\nu}=t_{\nu}; \beta) = \frac{e^{-\beta_{i}-\beta_{j}} \cdot \gamma_{t_{\nu}-2}^{(i,j)}}{\gamma_{t_{\nu}}} =: p_{\nu i, \nu j \mid t_{\nu}} \text{ for } i \neq j = 1, \dots, k.$$

So, the elements of the submatrix $I_g^{\beta\beta^r}$ of (28) are given by:

$$\left(\mathbf{I}_{g}^{\beta\beta^{\tau}}\right)_{ii} = \mathscr{C}\sum_{\nu=1}^{n} \left\{-p_{\nu i \mid l_{\nu}}\left(1-p_{\nu i \mid l_{\nu}}\right)+p_{\nu i}\left(1-p_{\nu i}\right)\right\}, \ i=1,\ldots,k \text{ and}$$
(29a)

$$\left(\mathbf{I}_{g}^{\beta\beta^{\mathsf{r}}}\right)_{ij} = \mathscr{E}\sum_{\nu=1}^{n} \left\{-p_{\nu i,\nu j \mid t_{\nu}} + p_{\nu i \mid t_{\nu}} \cdot p_{\nu j \mid t_{\nu}}\right\}, \ i \neq j = 1, \dots, k.$$
(29b)

The Fisher information matrix of g with respect to β and θ , as specified in (29ab) and (26cde), is not positive definite, but $I_g^{\theta\theta^r}$ is positive definite, and (14) can be used again to find the F-information in g with respect to β . This matrix is for i = 1, ..., k given by:

$$\left(I_{g}(\beta;\omega)\right)_{ii} = \mathscr{E}\sum_{\nu=1}^{n} \left\{-p_{\nu i \mid t_{\nu}}\left(1-p_{\nu i \mid t_{\nu}}\right)+p_{\nu i}\left(1-p_{\nu i}\right)\right\} - \sum_{\nu=1}^{n} \frac{\left[p_{\nu i}\left(1-p_{\nu i}\right)\right]^{2}}{\sum_{\ell=1}^{k} p_{\nu \ell}\left(1-p_{\nu \ell}\right)}, \quad (30a)$$

and for $i \neq j = 1, ..., k$

$$\left(\mathbf{I}_{g}(\beta;\omega)\right)_{ij} = \mathscr{E}\sum_{\nu=1}^{n} \left\{-p_{\nu i,\nu j \mid t_{\nu}} + p_{\nu i \mid t_{\nu}} \cdot p_{\nu j \mid t_{\nu}}\right\} - \sum_{\nu=1}^{n} \frac{p_{\nu i}(1-p_{\nu i}) \cdot p_{\nu j}(1-p_{\nu j})}{\sum_{\ell=1}^{k} p_{\nu \ell}(1-p_{\nu \ell})}.$$
 (30b)

Using (15ab) on the additivity of the F-information, the F-information in the conditional distribution $f(x \mid t; \beta)$ simply is found by subtracting the expressions in (30ab) from (27ab) and this gives:

$$\left(\mathscr{E}\operatorname{I}_{f}(\beta \mid T)\right)_{ii} = \mathscr{E}\sum_{\nu=1}^{n} p_{\nu i \mid t_{\nu}}\left(1 - p_{\nu i \mid t_{\nu}}\right) \text{ for } i = 1, \dots, k \text{ and}$$
(31a)

$$\left(\mathscr{E}\operatorname{I}_{f}(\beta \mid T)\right)_{ij} = \mathscr{E}\sum_{\nu=1}^{n} \left(p_{\nu i,\nu j \mid t_{\nu}} - p_{\nu i \mid t_{\nu}} \cdot p_{\nu j \mid t_{\nu}} \right) \text{ for } i \neq j = 1, \dots, k.$$
(31b)

All three relevant F-information matrices in the Rasch model are now specified in (27ab), (30ab) and (31ab). Once the expectations in (30ab) and (31ab) over the distribution of T_{ν} (22) are determined, they are easily computed for a given set of item and ability parameters. The diagonal of the F-information in the conditional distribution (31a), for instance, is

$$\left(\mathscr{E} \mathbf{I}_{f}(\beta \mid T)\right)_{ii} = \sum_{\nu=1}^{n} \sum_{t_{\nu}=0}^{k} p_{\nu i \mid t_{\nu}} \left(1 - p_{\nu i \mid t_{\nu}}\right) \exp(\theta_{\nu} t_{\nu}) \gamma_{t_{\nu}} \prod_{i=1}^{k} \left(1 - p_{\nu i}\right) \text{ for } i = 1, \dots, k.$$

The expressions (30ab) and (31b) change similarly.

Note that the expressions for the F-information matrices are all sums of independent contributions of the persons v = 1, ..., n.

Example

Consider the Rasch model with 3 items, $\beta_1 = -.5$, $\beta_2 = 0.0$ and $\beta_3 = 1.75$, and 4 persons with θ_v respectively -1.0, 0.0, 1.0, and 2.0.

The Fisher information matrix $(k + n \times k + n)$ is given by:

$$I_{p}(\omega) = \begin{cases} .689 & .000 & .000 & -.235 & -.235 & -.149 & -.070 \\ .000 & .748 & .000 & -.197 & -.250 & -.197 & -.105 \\ .000 & .000 & .647 & -.056 & -.126 & -.218 & -.246 \\ -.235 & -.197 & -.056 & .488 & .000 & .000 & .000 \\ -.235 & -.250 & -.126 & .000 & .611 & .000 & .000 \\ -.149 & -.197 & -.218 & .000 & .000 & .564 & .000 \\ -.070 & -.105 & -.246 & .000 & .000 & .421 \end{cases}$$

The structure of the matrix is clear: The item parameter submatrix as well as the ability parameter submatrix is diagonal and all coformations between items and abilities are negative.

The symmetric F-Information matrix in $p(k \times k)$, using (27ab), is given by

$$I_{p}(\beta;\omega) = \begin{pmatrix} .435 & -.260 & -.174 \\ -.260 & .472 & -.212 \\ -.174 & -.212 & .386 \end{pmatrix}.$$

Note that in JML estimation in the Rasch model the diagonal of the item parameter submatrix of the Fisher information matrix is sometimes used for estimating the standard errors of the item parameters (Wright, 1977). (In the computations the matrix is evaluated at the estimates). This practice leads to an underestimate of the standard error, because a better estimate, taking account of the negative coformation between the items and abilities, would use the diagonal of the F-information matrix.

The F-information matrices in the conditional distribution f and the marginal distribution g are given by

$$\mathbf{I}_{f}(\beta;\omega) = \begin{bmatrix} .409 & -.272 & -.137 \\ -.272 & .451 & -.180 \\ -.137 & -.180 & .316 \end{bmatrix} \text{ and } \mathbf{I}_{g}(\beta;\omega) = \begin{bmatrix} .026 & .012 & -.037 \\ .012 & .021 & -.032 \\ -.037 & -.032 & .070 \end{bmatrix}.$$

The loss of information in using the conditional distribution $f(x \mid t;\beta)$ instead of the full model $p(x;\beta,\theta)$, that is, in using CML instead of JML estimation for the item parameters is given by $I_g(\beta;\omega)$, the F-information matrix in g. As a scalar measure of this loss, the trace of this matrix is used. Note that in our application it is easily checked that the determinant function is useless for this, because all F-information matrices are singular:

$$\sum_{i} I_{p}(\beta; \omega) = \sum_{i} I_{p}(\beta; \omega) = \sum_{i} I_{f}(\beta; \omega) = \sum_{i} I_{f}(\beta; \omega) = \sum_{i} I_{g}(\beta; \omega) = \sum_{i} I_{g}(\beta; \omega) = 0.$$

The loss of information in the example appears to be very small. In the example, the efficiency measure, defined in (17), $\text{EFF}(\beta;f:p) = \text{tr}(I_f(\beta;\omega))/\text{tr}(I_p(\beta;\omega))$, expressing the loss in using the conditional instead of the full model is $\text{EFF}(\beta;f:p) = 1.176/1.293 = .910.\square$

Efficiency comparison of JML and CML estimation

The efficiency of CML versus JML will be reported for some typical item parameter and ability parameters sets. For these sets, 100 ability parameters were drawn from a normal ability distribution $N(\mu, \sigma^2)$. Table 1 gives the efficiencies for 10 items with $\beta = (-3, -2, -1.5, -1, -0.5, 0.5, 1, 1.5, 2, 3)$ and 100 abilities drawn from the normal distribution with varying μ and σ^2 .

	σ^2			
μ	.25	1	2.25	4
-2	.984	.985	.985	.986
-1	.988	.987	.987	.986
0	.988	.988	.987	.986
1	.988	.988	.986	.986
2	.985	.985	.986	.985

Table 1

It should be noted that the traces of $I_{f}(\beta;\omega)$ as well as of $I_{p}(\beta;\omega)$ decrease, when σ^{2} is increased, and also when the distance between the mean of the ability- and the mean of the item parameters, $|\mu - \frac{1}{k} \Sigma_k \beta_i|$, is increased. However, the efficiency of CML versus JML, reported in Table 1, hardly shows any variation and amounts to at least 98.4%.

Table 2 shows the relation between the efficiency and the spread in the item parameters. For 100 abilities drawn from N(0,1), the efficiency for estimating 10 equidistant item parameters with $\sum_{i} \beta_{i} = 0$ and varying stepsize, $\beta_{i+1} - \beta_{i}$, i = 1, ..., 9is given.

Table 2 EFF $(\beta; f; p)$ for 10 equidistant items with $\sum_i \beta_i = 0$ and varying stepsize and 100 abilities from N(0,1)

step	$\mathrm{EFF}(\beta; f; p)$
1	.9451
.5	.9891
.25	.9968
.125	.9992
.0625	.9998
0	1.000

The efficiency of CML versus JML estimation increases with decreasing spread of the item parameters. In case all item parameters are equal there is no loss of information. The theoretical explanation for this is that the weak ancillarity condition in Theorem 3

(23) is fulfilled in this special case: $\partial \ln \prod_{\nu} \gamma_{t_{\nu}}(\beta) / \partial \beta_j = \sum_{\nu} t_{\nu} / k$ for j = 1, ..., k and $I_g(\beta; \omega) = 0$.

In Table 3 the relation between the efficiency and the test length is illustrated. For 100 abilities drawn from N(0.5,1), the efficiency is reported for estimating the item parameters $\beta = (0, 1, 2)$, which are *n* times in a test. The test length is equal to 3n.

Table 3		
EFF(β ; f: p) for $\beta = (0, 1, 2)$ with varying tes	t lengt	h
and 100 abilities from $N(0.5,1)$		

length	$EFF(\beta; f: p)$
3	.9360
6	.9909
9	.9965
12	.9982
15	.9989
18	.9992
21	.9995

An increasing efficiency is observed when the test length is increased. It is clear that already at a relative short test length of 21 items, CML estimation is almost as efficient as JML estimation of the item parameters.

F-Information in the Dichotomous Rasch Model: Comparing MML and CML

In order to compare CML with MML estimation of the item parameters in the Rasch model, the marginal likelihood function (3) in which the θ_v 's are not fixed ability parameters but random draws from the ability distribution $h(\theta; \xi)$, is considered. With again $T_v = \sum_i X_{vi}$, this can be rewritten as:

$$p_{m}(x;\beta,\xi) = \prod_{\nu} p_{m}(x_{\nu};\beta,\xi) = \prod_{\nu} \int_{-\infty}^{\infty} \prod_{i} P(X_{\nu i} = x_{\nu i} \mid \theta_{\nu};\beta_{i}) h(\theta_{\nu};\xi) d\theta_{\nu} = \prod_{\nu} f(x_{\nu} \mid t_{\nu};\beta) \prod_{\nu} \int_{-\infty}^{\infty} g(t_{\nu} \mid \theta_{\nu};\beta) h(\theta_{\nu};\xi) d\theta_{\nu}.$$
(32)

In (32) the parameters are the item parameter β and the parameter of the ability distribution ξ . This two-parameter distribution is factored in a part which is only dependent on the interest parameter β , and a part which is the distribution of T, which depends on both parameters β and ξ . For estimating the item parameters with CML, as seen in (5), only the first part of the factorization is used. In MML estimation the full model (32) is used. And for comparing CML estimation with MML estimation, the F-information in f, $\mathscr{E} I_f(\beta \mid T)$, given in (31ab), can be compared with the F-information in p_m with respect to β .

Note that because in the model (32) every person is a random draw from an ability distribution, it suffices to derive the F-information in p_m for only one person, indexed with v. The same is true for the F-information in f (see (31ab)).

In order to determine these F-information matrices completely, an ability distribution $h(\theta_{\nu}, \xi)$ has to be specified. Here it is assumed that the ability distribution is normal with parameter $\xi = (\mu, \sigma^2)$, and the density is given by

$$h(\theta_{\nu};\xi) = \sigma^{-1} \phi((\theta_{\nu} - \mu) / \sigma), \qquad (33)$$

with $\phi(y)$ being the standard normal density:

$$\phi(y) = (2\pi)^{-\frac{1}{2}} \exp(-y^2/2).$$
(34)

F-information in p_m

As before, this matrix can be determined once the Fisher information matrix has been obtained. If the length of the vector parameter of the ability distribution is ℓ , this matrix can be written in the following partitioned form:

$$\mathbf{I}_{p_{m}}(\omega) = \mathscr{C} \begin{bmatrix} S_{p;\beta} S_{p;\beta}^{T} & S_{p;\beta} S_{p;\xi}^{T} \\ S_{p;\xi} S_{p;\beta}^{T} & S_{p;\xi} S_{p;\xi}^{T} \end{bmatrix} =: \begin{bmatrix} \mathbf{I}_{p_{m}}^{\beta\beta^{T}} & \mathbf{I}_{p_{m}}^{\beta\xi^{T}} \\ \mathbf{I}_{p_{m}}^{\xi\beta^{T}} & \mathbf{I}_{p_{m}}^{\xi\xi^{T}} \end{bmatrix}$$
(35)

The submatrices, $I_{p_m}^{\beta\beta^r}(k \times k)$, $I_{p_m}^{\beta\xi^r} = [I_{p_m}^{\xi\beta^r}]^T(k \times \ell)$ and $I_{p_m}^{\xi\xi^r}(\ell \times \ell)$ are again the negative expectation of the second derivatives of $\ln p_m$ with respect to the parameters. For notational convenience define:

$$Q_{\nu} := Q_{\nu} (\theta_{\nu}, x_{\nu i}, t_{\nu}, \beta, \mu, \sigma^{2}) = \exp(\theta_{\nu} t_{\nu}) \prod_{i=1}^{k} (1 - p_{\nu i}) h(\theta_{\nu}; \xi),$$
(36)

in the normal density ((33) with (34)), can be substituted for $h(\theta_{\nu}; \xi)$. So, considered is the model:

$$P_{m}(x_{\nu};\beta,\xi) = f(x_{\nu} \mid t_{\nu};\beta) \cdot \int_{-\infty}^{\infty} g(t_{\nu} \mid \theta_{\nu};\beta) h(\theta_{\nu};\xi) d\theta_{\nu}, \qquad (37)$$

in which (see (22), with (25) and (36)),

$$g(t_{\nu} \mid \theta_{\nu}; \beta) = \exp(\theta_{\nu} t_{\nu}) \cdot \gamma_{t_{\nu}}(\beta) \cdot \prod_{i=1}^{k} (1 - p_{\nu i}) = Q_{\nu} \gamma_{t_{\nu}} / h(\theta_{\nu}; \xi), \qquad (38)$$

and (divide (21) by (22) for one person)

$$f(x_{v} \mid t_{v}; \beta) = \frac{\exp\left(-\sum_{i} \beta_{i} x_{vi}\right)}{\gamma_{t_{v}}(\beta)}.$$
(39)

The marginal distribution of T_{ν} is given by

$$P(T_{\nu=t_{\nu}};\beta,\xi) = \int_{-\infty}^{\infty} g(t_{\nu} \mid \theta_{\nu};\beta) h(\theta_{\nu};\xi) d\theta_{\nu} = \gamma_{t_{\nu}} \int Q_{\nu} d\theta_{\nu}.$$
(40)

Inserting (38) and (39) in (37) and taking the logarithm, gives

$$\ln p_m(x_{\nu};\beta,\xi) = \sum_i \beta_i x_{\nu i} + \ln \int_{-\infty}^{\infty} \exp(\theta_{\nu} t_{\nu}) h(\theta_{\nu};\xi) \prod_{i=1}^{k} (1-p_{\nu i}) d\theta_{\nu} = \sum_i \beta_i x_{\nu i} + \ln \int Q_{\nu} d\theta_{\nu}.$$

With some algebra the following expressions are easily derived:

$$\frac{\partial Q_{\nu}}{\partial \beta_j} = Q_{\nu} p_{\nu j}, \text{ for } j = 1, \dots, k,$$

$$\begin{aligned} \frac{\partial Q_{\nu}}{\partial \mu} &= Q_{\nu}(\theta_{\nu} - \mu)\sigma^{-2}, \\ \frac{\partial Q_{\nu}}{\partial \sigma^{2}} &= Q_{\nu} \left[-\frac{1}{2}\sigma^{-2} + \frac{1}{2}(\theta_{\nu} - \mu)^{2}\sigma^{-4} \right], \\ \frac{\partial^{2} Q_{\nu}}{\partial \beta_{j}^{2}} &= Q_{\nu}(2p_{\nu j} - 1)p_{\nu j}, \text{ for } j = 1, ..., k, \end{aligned}$$
(41a)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial \beta_{j}^{2}\partial \beta_{m}} &= Q_{\nu}p_{\nu j}p_{\nu m}, \text{ for } m \neq j = 1, ..., k, \end{aligned}$$
(41b)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial \mu^{2}} &= Q_{\nu} \left[(\theta_{\nu} - \mu)^{2} \sigma^{-4} - \sigma^{-2} \right], \end{aligned}$$
(42a)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial (\sigma^{2})^{2}} &= Q_{\nu} \left[\frac{3}{4} \sigma^{-4} - \frac{3}{2}(\theta_{\nu} - \mu)^{2} \sigma^{-6} + \frac{1}{4}(\theta_{\nu} - \mu)^{4} \sigma^{-8} \right], \end{aligned}$$
(42b)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial \sigma^{2} \partial \mu} &= Q_{\nu} \left[-\frac{3}{2}(\theta_{\nu} - \mu) \sigma^{-4} + \frac{1}{2}(\theta_{\nu} - \mu)^{3} \sigma^{-6} \right], \end{aligned}$$
(42c)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial \mu^{2} \partial \beta_{j}} &= Q_{\nu} p_{\nu j} (\theta_{\nu} - \mu) \sigma^{-2}, \text{ for } j = 1, ..., k, \end{aligned}$$
(43a)
$$\begin{aligned} \frac{\partial^{2} Q_{\nu}}{\partial \sigma^{2} \partial \beta_{j}} &= Q_{\nu} p_{\nu j} \left[-\frac{1}{2} \sigma^{-2} + \frac{1}{2}(\theta_{\nu} - \mu)^{2} \sigma^{-4} \right] \text{ for } j = 1, ..., k. \end{aligned}$$
(43b)

The second derivatives of $\ln p_m(x_v;\beta,\xi)$ are generally written as (see Appendix):

$$\frac{\partial^2 \ln p_m(x_\nu;\beta,\xi)}{\partial a \partial b} = \frac{\int \frac{\partial^2 Q_\nu}{\partial a \partial b} d\theta_\nu}{\int Q_\nu d\theta_\nu} - \frac{\int \frac{\partial Q_\nu}{\partial a} d\theta_\nu \int \frac{\partial Q_\nu}{\partial b} d\theta_\nu}{\left(\int Q_\nu d\theta_\nu\right)^2},\tag{44}$$

with a and b being any pair from the parameters μ , σ^2 , and β_j , j = 1, ..., k. Taking minus the expectation of (44) over the distribution of T_v (40) gives:

$$-\mathscr{E}\frac{\partial^{2}\ln p_{m}(x_{\nu};\beta,\xi)}{\partial a\,\partial b} = -\sum_{t_{\nu}=0}^{k}\gamma_{t_{\nu}}\left[\int\frac{\partial^{2}Q_{\nu}}{\partial a\,\partial b}d\theta_{\nu} - \frac{\int\frac{\partial Q_{\nu}}{\partial a}d\theta_{\nu}\int\frac{\partial Q_{\nu}}{\partial b}d\theta_{\nu}}{\int Q_{\nu}d\theta_{\nu}}\right],\qquad(45)$$

which is the general expression for all parts of the Fisher information matrix (35). Using (41ab) in (45) gives the item parameter part:

$$\left(\mathbf{I}_{p_{m}}^{\beta\beta^{T}}\right)_{jj} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left[\int (2p_{\nu j} - 1)p_{\nu j}Q_{\nu}d\theta_{\nu} - \frac{\left(\int p_{\nu j}Q_{\nu}d\theta_{\nu}\right)^{2}}{\int Q_{\nu}d\theta_{\nu}} \right]$$
(46a)

for j = 1, ..., k, and

$$\left(\mathbf{I}_{p_{m}}^{\beta\beta^{T}}\right)_{mj} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left[\int p_{\nu m} p_{\nu j} Q_{\nu} d\theta_{\nu} - \frac{\int p_{\nu m} Q_{\nu} d\theta_{\nu} \int p_{\nu j} Q_{\nu} d\theta_{\nu}}{\int Q_{\nu} d\theta_{\nu}}\right]$$
(46b)

for $m \neq j = 1, \dots, k$.

Using (42abc) in (45) gives the ability distribution parameter part:

$$\left(\mathbf{I}_{p_{m}}^{\xi\xi^{T}}\right)_{11} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left[\int \mathcal{Q}_{\nu} \left[\left(\theta_{\nu} - \mu\right)^{2} \sigma^{-4} - \sigma^{-2} \right] d\theta_{\nu} - \frac{\left(\int \mathcal{Q}_{\nu} \left(\theta_{\nu} - \mu\right) \sigma^{-2} d\theta_{\nu}\right)^{2}}{\int \mathcal{Q}_{\nu} d\theta_{\nu}} \right]$$
(47a)

$$\left(\mathbf{I}_{p_{m}}^{\xi\xi^{T}} \right)_{22} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left\{ \int \mathcal{Q}_{\nu} \left[\frac{3}{4} \, \sigma^{-4} - \frac{3}{2} \left(\theta_{\nu} - \mu \right)^{2} \sigma^{-6} + \frac{1}{4} \left(\theta_{\nu} - \mu \right)^{4} \sigma^{-8} \right] d\theta_{\nu} - \frac{\left[\int \mathcal{Q}_{\nu} \left[-\frac{1}{2} \, \sigma^{-2} + \frac{1}{2} \left(\theta_{\nu} - \mu \right)^{2} \, \sigma^{-4} \right] d\theta_{\nu} \right]^{2} \right] \right\}$$

$$\left(\mathbf{I}_{p_{m}}^{\xi\xi^{T}} \right)_{12} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left\{ \int \mathcal{Q}_{\nu} \left[-\frac{3}{2} \left(\theta_{\nu} - \mu \right) \sigma^{-4} + \frac{1}{2} \left(\theta_{\nu} - \mu \right)^{3} \sigma^{-6} \right] d\theta_{\nu} - \frac{1}{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \right)^{2} \right] d\theta_{\nu} - \frac{1}{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \right)^{2} \right) d\theta_{\nu} - \frac{1}{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \right)^{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \right)^{2} \left(\theta_{\nu} - \theta_{\nu} \right)^{2} \left(\theta_{\nu} -$$

$$\frac{\int Q_{\nu}(\theta_{\nu}-\mu)\sigma^{-2}d\theta_{\nu} \int Q_{\nu}\left[-\frac{1}{2}\sigma^{-2}+\frac{1}{2}(\theta_{\nu}-\mu)^{2}\sigma^{-4}\right]d\theta_{\nu}}{\int Q_{\nu}d\theta_{\nu}} \right\}.$$
(47c)

Finally, using (43ab) in (45) gives the elements of $I_{p_m}^{\beta\xi^r}$. For j=1,...,k

$$\left(I_{p_{m}}^{\beta\xi^{T}}\right)_{1j} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left\{ \int Q_{\nu} P_{\nu j} (\theta_{\nu} - \mu) \sigma^{-2} d\theta_{\nu} - \frac{\int Q_{\nu} P_{\nu j} d\theta_{\nu} \int Q_{\nu} (\theta_{\nu} - \mu) \sigma^{-2} d\theta_{\nu}}{\int Q_{\nu} d\theta_{\nu}} \right\}, \quad (48a)$$

and

$$\left(\mathbf{I}_{\rho_{m}}^{\beta\xi^{\tau}} \right)_{2j} = -\sum_{t_{\nu}=0}^{k} \gamma_{t_{\nu}} \left\{ \int \mathcal{Q}_{\nu} p_{\nu j} \left[-\frac{1}{2} \sigma^{-2} + \frac{1}{2} (\theta_{\nu} - \mu)^{2} \sigma^{-4} \right] d\theta_{\nu} - \frac{\int \mathcal{Q}_{\nu} p_{\nu j} d\theta_{\nu} \int \mathcal{Q}_{\nu} \left[-\frac{1}{2} \sigma^{-2} + \frac{1}{2} (\theta_{\nu} - \mu)^{2} \sigma^{-4} \right] d\theta_{\nu} - \frac{\int \mathcal{Q}_{\nu} p_{\nu j} d\theta_{\nu} \int \mathcal{Q}_{\nu} d\theta_{\nu} \right\}$$

$$(48b)$$

With respectively (46ab), (47abc), and (48ab), all expressions for the submatrices of the Fisher information matrix (35) are obtained. From these expressions the F-information matrix $I_{p_m}(\beta;\omega)$ can be computed. Because $I_{p_m}^{\xi\xi^{T}}$ is positive definite, the F-information matrix with respect to β is given by (14), that is,

$$I_{p_{m}}(\beta;\omega) = I_{p_{m}}^{\beta\beta^{T}} - I_{p_{m}}^{\beta\xi^{T}} [I_{p_{m}}^{\xi\xi^{T}}]^{-1} I_{p_{m}}^{\xi\beta^{T}} .$$

F-information in f(x|t)

In order to make the relevant comparisons, the expressions for the F-information in the conditional distribution are required. Taking expectations over the distribution of T_{ν} (40) in (31ab) gives:

$$\left(\mathbf{I}_{f}(\boldsymbol{\beta};\boldsymbol{\omega})\right)_{jj} = \left(\mathscr{E}\mathbf{I}_{f}(\boldsymbol{\beta} \mid T)\right)_{jj} = \sum_{t_{v}=0}^{k} \left\{ e^{-\beta_{j}} \cdot \gamma_{t_{v}-1}^{(j)} - \frac{\left[e^{-\beta_{j}} \cdot \gamma_{t_{v}-1}^{(j)}\right]^{2}}{\gamma_{t_{v}}} \right\} \int \mathcal{Q}_{v} d\theta_{v}$$

for j = 1, ..., k, and

$$\left(\mathbf{I}_{f}(\boldsymbol{\beta};\boldsymbol{\omega})\right)_{mj} = \left(\mathscr{E}\mathbf{I}_{f}(\boldsymbol{\beta} \mid T)\right)_{mj} = \sum_{t_{\nu}=0}^{k} \left\{-\frac{e^{-\beta_{m}-\beta_{j}} \cdot \gamma_{t_{\nu}-1}^{(m)} \gamma_{t_{\nu}-1}^{(j)}}{\gamma_{t_{\nu}}} + e^{-\beta_{m}-\beta_{j}} \cdot \gamma_{t_{\nu}-2}^{(m,j)}\right\} \int Q_{\nu} d\theta_{\nu}$$

for $m \neq j = 1, \dots, k$.

Example

Consider the example similar to the one used in the computations of the information matrices in the Rasch model with fixed ability parameters. For the Rasch model with 3 items, $\beta_1 = -.5$, $\beta_2 = 0.0$ and $\beta_3 = 1.75$, and assuming that θ is distributed as $N(\mu, \sigma^2)$, with $\mu = 0$ and $\sigma^2 = 1$, the expected Fisher information matrix, $(k + \ell) \times (k + \ell)$ (35), for one random draw of the ability distribution is given by:

$$\mathbf{I}_{p_{m}}(\omega) = \begin{pmatrix} \mathbf{I}_{p_{m}}^{\beta\beta^{T}} & \mathbf{I}_{p_{m}}^{\beta\xi^{T}} \\ \mathbf{I}_{p_{m}}^{\xi\beta^{T}} & \mathbf{I}_{p_{m}}^{\xi\xi^{T}} \end{pmatrix} = \begin{pmatrix} .172 & -.027 & -.016 & -.129 & .016 \\ -.027 & .178 & -.018 & -.134 & .002 \\ -.016 & -.018 & .119 & -.085 & -.029 \\ -.129 & -.134 & -.085 & .348 & .011 \\ .016 & .002 & -.029 & .011 & .038 \end{pmatrix}$$

As expected, the upper left 3×3 submatrix $I_{p_m}^{\beta\beta^r}$ is symmetric, and the coformation elements between the items are negative. This matrix is the information matrix with respect to β , when μ and σ^2 are considered known. In the submatrix $I_{p_m}^{\beta\xi^r}$, negative coformations between the item parameters and the mean of the ability distribution are obtained, while the coformation between items and the variance of the ability

distribution has no constant sign. In the part of the parameters of the ability distribution, a non-zero coformation between μ and σ^2 can be observed. Note that only in case the mean of the item parameters equals the mean of the ability distribution, the coformation between μ and σ^2 is zero.

The F-information matrices in p_m and in the conditional distribution f are given by

$$\mathbf{I}_{p_m}(\beta;\omega) = \begin{pmatrix} .114 & -.080 & -.034 \\ -.080 & .126 & -.046 \\ -.034 & -.046 & .080 \end{pmatrix} \mathbf{I}_f(\beta;\omega) = \mathscr{E}I_f(\beta \mid T) = \begin{pmatrix} .114 & -.080 & -.034 \\ -.080 & .122 & -.042 \\ -.034 & -.042 & .076 \end{pmatrix}.$$

The efficiency for estimating the item parameters β using the conditional model (CML) instead of the full model (MML) can be computed by dividing the traces of the F-information matrices (17). In this example EFF(β ; $f:p_m$) = .3115/.3197 = .9744.

Efficiency comparison CML and MML estimation

The efficiency of CML versus MML will be reported for some typical item parameter sets. In Table 4, the traces of the F-information matrices and the efficiency for 10 items with $\beta = (-3, -2, -1.5, -1, -0.5, 0.5, 1, 1.5, 2, 3)$ with varying μ and σ^2 are presented.

It can be seen that the traces of the F-information matrices, and, of course, also the efficiencies, are symmetric in μ . Furthermore it is seen, that, as expected, the traces of $I_f(\beta; \omega)$ as well as $I_{p_m}(\beta; \omega)$ decrease, σ^2 is increased, also when the distance between the mean of the ability- and the mean of the item parameters, $|\mu - \Sigma_k \beta_i / k|$, is increased.

Table 4

	σ^2			
μ	.25	1	2.25	4
$\operatorname{tr}(\mathbf{I}_{\mathbf{f}}(\boldsymbol{\beta};\boldsymbol{\omega}))$				
-2	.976	.939	.894	.843
-1	1.188	1.136	1.061	.976
0	1.255	1.204	1.122	1.024
1	1.888	1.136	1.061	.976
2	.976	.939	.894	.843
		$\operatorname{tr}(\mathbf{I}_{p_{-}}(\boldsymbol{\beta};\boldsymbol{\omega}))$	v))	
-2	.980	.947	.903	.854
-1	1.193	1.144	1.073	.989
0	1.258	1.212	1.134	1.039
1	1.193	1.144	1.073	.989
2	.980	.947	.903	.854
		EFF ($\beta; f: \mu$	<i>p</i> _m)	
-2	.995	.992	.989	.987
-1	.997	.993	.989	.987
0	.997	.993	.989	.989
1	.997	.993	.989	.987
2	.995	.992	.989	.987

tr(I_f($\beta; \omega$)), tr(I_{p_m}($\beta; \omega$)) and EFF($\beta; f: p_m$) for $\beta = (-3, -2, -1.5, -1, -0.5, 0.5, 1, 1.5, 2, 3)$ and ability from $N(\mu, \sigma^2)$

The efficiency of CML versus MML decreases somewhat with increasing σ^2 , and decreasing $|\mu - \Sigma_k \beta_i / k|$. However, the efficiency of CML versus MML amounts to at least 98.7%.

In Table 5, the relation between the efficiency and the spread in the item parameters is given. For 1 ability drawn from N(0,1), the efficiency for estimating 10 equidistant item parameters with $\sum_{i} \beta_{i} = 0$ and varying stepsize, $\beta_{i+1} - \beta_{i}$, i = 1, ...9 is given.

Table	5
-------	---

EFF(β ; f:p_m) for $\Sigma_i \beta_i = 0$ with varying stepsize and 1 ability from N(0,1)

step	$EFF(\beta; f: p_m)$
1	.9585
.5	.9945
.25	.9994
.125	.9999
.0625	1.0000

Although there is hardly any loss in information, the efficiency of CML versus MML estimation increases with a decrease in the spread between the item parameters.

In Table 6, the relation between the efficiency and the test length is illustrated. For 1 ability drawn from N(0.5,1), the efficiency is reported for estimating the item parameters $\beta = (0,1,2)$, which are *n* times in a test. The test length is 3n.

Table 6
EFF(β ; f:p _m) for $\beta_i = 0, 1, 2$ with varying test length
and 1 ability from $N(0.5,1)$

length	$EFF(\beta; f: p_m)$
3	.9822
6	.9983
9	.9993
12	.9996
15	.9998
18	1.0000

An increasing efficiency is observed when the test length is increased. Already with a relative short test of 18 items, CML estimation is as efficient as MML estimation of the item parameters.

Conclusion

In this paper it is shown that the concept of F-information, a generalization of Fisher information, is a useful tool for evaluating the loss of information in conditional maximum likelihood estimation. In separable two-parameter models with a sufficient statistic for the nuisance parameter, the conditional model which only depends on the interest parameter is often used instead of the full two-parameter model for estimating the interest parameter. With the F-information concept it is possible to investigate the conditions under which there is no loss of information in CML estimation, and furthermore, if there is a loss, to quantify this.

In this paper the main properties of F-information are presented and the conditions for no loss of information are specified. Especially in the case of exponential family models, it is shown that these conditions can be easily checked. It is shown that in the Poisson Counts Model these conditions are met. For the Rasch model for dichotomously scored items these conditions are not fulfilled, which means that loss of information in CML estimation of the item parameters in this model is to be expected.

For the dichotomous Rasch model the expressions needed to be able to investigate the loss of information using CML estimation of the item parameters are derived in detail. For the Rasch model with fixed ability parameters, a comparison was made between JML and CML estimation, and, under the assumption of a normal ability distribution, a comparison of MML and CML estimation. The comparisons were made in several conditions. Varied were: the spread of the item difficulty parameters, the mean and the variance of the ability distribution and the test length. In almost all comparisons, some loss of information in using CML appeared. However, in all the comparisons of CML to JML as well as to MML the loss was very small. The reported efficiencies are always larger than 93%, and for tests with 20 or more items larger than 99%.

On basis of these results it can be concluded that CML item parameter estimation in the Rasch model, which has some other known practical and theoretical attractive properties, is also from an information point of view a sound practice. Hardly any loss of information is to be expected compared to alternative estimation methods.

The method of efficiency comparison described in the paper, applies to any separable two-parameter model. For popular extensions of the dichotomous Rasch model (Fischer & Molenaar, 1995), for example the partial credit model and the one-parameter logistic model, comparable work is in progress.

References

- Andersen, E.B. (1970). Asymptotic properties of conditional maximum likelihood estimators. Journal of the Royal Statistical Society, Series B, 32, 283-301
- Andersen, E.B. (1973). Conditional inference and models for measuring. Copenhagen: Mentalhygiejnisk Forlag.
- Basu, D. (1977). On the elimination of nuisance parameters. Journal of the American Statistical Society, 72, 355-366.
- Bhapkar, V.P. (1989). Conditioning on ancillary statistics and loss of information in the presence of nuisance parameters. *Journal of Statistical Planning and Inference*, 21, 139-160.
- Bhapkar, V.P. (1991). Loss of information in the presence of nuisance parameters and partial sufficiency. *Journal of Statistical Planning and Inference*, 28, 185-203.
- Efron, B. (1977). On the efficiency of Cox's likelihood function for censored data. Journal of the American Statistical Society, 72, 557-565.
- Engelen, R.J.H. (1989). *Parameter estimation in the logistic item response model*. (Doctoral Thesis.) Enschede: University of Twente.
- Fischer, G.H. (1974). *Einführung in die Theorie psychologischer Tests*. [Introduction to mental test theory.] Bern: Huber.
- Fischer, G.H., & Molenaar, I.W. (Eds.) (1995). Rasch Models. New York: Springer.
- Glas, C.A.W. (1989). Contributions to estimating and testing Rasch models. (Doctoral Thesis.) Enschede: University of Twente.
- Liang, K. (1983). On information and ancillary in the presence of a nuisance parameter. *Biometrika*, 70, 607-612.
- Louis, T.A. (1982). Finding the observed information matrix with the EM algorithm. Journal of the Royal Statistical Society, Series B, 44, 226-233.
- Molenaar, I.W. (1995). Estimation of item parameters. In: Fischer, G.H., & Molenaar, I.W. (Eds.). Rasch Models (pp.39-51). New York: Springer.

Rao, C.R. (1973). Linear statistical inference and its applications. New York: Wiley.

Rasch, G. (1960). Probabilistic models for some intelligence and attainment tests.
Copenhagen: The Danish Institute of Educational Research. (Expanded edition, 1980. Chicago: The University of Chicago Press.)

- Wright, B.J. (1977). Solving measurement problems with the Rasch model. Journal of Educational Measurement, 14, 97-116.
- Zwinderman, A.H. (1991). Studies of estimating and testing Rasch models. (NICI Technical Report 91-01.) (Doctoral Thesis.) Nijmegen: University of Nijmegen.

40

*

Appendix

If the ability distribution $h(\theta_{\nu}; \xi)$ is considered as a prior distribution, then the a posteriori distribution of the ability of person ν , θ_{ν} , given the score t_{ν} , is by Bayes rule equal to:

$$k(\theta_{\nu} \mid t_{\nu}; \beta, \xi) =: \frac{g(t_{\nu} \mid \theta_{\nu}; \beta) h(\theta_{\nu}; \xi)}{\int_{-\infty}^{\infty} g(t_{\nu} \mid \theta_{\nu}; \beta) h(\theta_{\nu}, \xi) d\theta_{\nu}} = \frac{Q_{\nu}}{\int Q_{\nu} d\theta_{\nu}}.$$

The second derivatives with respect to the item parameters can be rewritten as:

$$\frac{\partial^2 \ln p_m(x_v;\beta,\xi)}{\partial \beta_j^2} = \int_{-\infty}^{\infty} (2p_{vj} - 1) p_{vj} k(\theta_v \mid t_v;\beta,\xi) d\theta_v - \left[\int_{-\infty}^{\infty} p_{vj} k(\theta_v \mid t_v;\beta,\xi) d\theta_v \right]^2$$
(i)

for j = 1, ..., k.

and

$$\frac{\partial^{2} \ln p_{m}(x_{\nu};\beta,\xi)}{\partial \beta_{m} \partial \beta_{j}} = \int_{-\infty}^{\infty} p_{\nu m} p_{\nu j} k(\theta_{\nu} \mid t_{\nu};\beta,\xi) d\theta_{\nu} - \left[\int_{-\infty}^{\infty} p_{\nu m} k(\theta_{\nu} \mid t_{\nu};\beta,\xi) d\theta_{\nu}\right] \left[\int_{-\infty}^{\infty} p_{\nu j} k(\theta_{\nu} \mid t_{\nu};\beta,\xi) d\theta_{\nu}\right]$$
(ii)

for $m \neq j = 1, \dots, k$.

(i) and (ii) can be written as moments in the a posteriori distribution of θ_{ν} given t_{ν} .

Because $(2p_{\nu j} - 1)p_{\nu j} = -(1 - p_{\nu j})p_{\nu j} + p_{\nu j}^{2}$, it yields that:

$$\frac{\partial^2 \ln p_m(x_{\nu};\beta,\xi)}{\partial \beta_j^2} = - \mathscr{E}\left(p_{\nu j}(1-p_{\nu j}) \mid t_{\nu}\right) + \operatorname{Var}\left(p_{\nu j} \mid t_{\nu}\right)$$

$$\frac{\partial^2 \ln p_m(x_{\nu};\beta,\xi)}{\partial \beta_m \partial \beta_j} = \operatorname{Cov}(p_{\nu m}, p_{\nu j} \mid t_{\nu}).$$

Evaluating the negatives of these expressions in the estimates of the parameters, the item parameter part of the observed information matrix of the Rasch model using MML estimation is obtained. The same expressions were deduced by Glas (1989, p. 50-55, Appendix D), using the theory of Louis (1982) on the relation between the observed and expected information matrix

Recent Measurement and Research Department Reports:

- 97-1 H.H.F.M. Verstralen. A Logistic Latent Class Model for Multiple Choice Items.
- 97-2 N.D. Verhelst. A New Heuristic for Estimating the Reliability of a Test.
- 97-3 N.D. Verhelst and R.J.H. Engelen. Estimating Latent Abilities from Raw Test Scores.
- 97-4 G.J.J.M. Straetmans and T.J.H.M. Eggen. Comparison of Test Administration Procedures for Placement Decisions in a Mathematics Course.
- 97-5 H.H.F.M. Verstralen. A Latent IRT Model for Options of Multiple Choice Items.