

DIF Detection via Regularization

Ruoyi Zhu

1/6/2020

Contents

1	DIF Detection Methods	1
1.1	Likelihood Ratio Test	1
2	Multidimensional 2PL Model with DIF	1
3	Graded Response Model with DIF	2
4	Model Identifiability Constraint	5
5	Uniform DIF Detection via LASSO	6
5.1	E step	6
5.2	M step	8
5.3	Simulation	13
6	Non-uniform DIF Detection via LASSO	16
6.1	E step	16
6.2	M step	16
6.3	Simulation	18
7	Non-uniform DIF Detection via Group LASSO	22
7.1	E step	22
7.2	M step	22
7.3	Simulation	25

1 DIF Detection Methods

1.1 Likelihood Ratio Test

In likelihood ratio test, DIF will be tested one item at a time. The null hypothesis is that all but the anchor items are DIF. To test whether the j th item is DIF, the alternative hypothesis is that the anchor items and the j th item is DIF free, and all other items are DIF. If the null hypothesis is rejected (adjust p-value < 0.05), the tested item will be flagged as DIF. Since we are testing several hypotheses, multiple comparisons need to be performed. The adjust p-values are the family-wise error rate. We assume the latent traits in reference and focal groups follow normal distributions. The mean and variance of reference group are fixed to be 0 and 1, respectively. The means and variances of focal groups can be freely estimated.

2 Multidimensional 2PL Model with DIF

Assume a total sample size $N = N_1 + N_2 + N_3$, where N_1 , N_2 and N_3 are sample sizes of the reference group and two focal groups, respectively. Also assume a test length m and trait dimension q . For a dichotomously scored item j , the probability that examinee i with ability vector θ_i giving a correct response to item j is

$$P_j(\boldsymbol{\theta}_i) = \frac{1}{1 + e^{-(\mathbf{a}_j^T \boldsymbol{\theta}_i + d_j + (\mathbf{y}_i \boldsymbol{\gamma}_j) \boldsymbol{\theta}_i + \mathbf{y}_i \boldsymbol{\beta}_j)}} (i = 1, \dots, N; j = 1, 2, \dots, m). \quad (1)$$

\mathbf{y}_i is a group indicator including all the grouping information related to DIF. $\mathbf{y}_i = (0, 0)$ if examinee i is in the reference group, $\mathbf{y}_i = (1, 0)$ if examinee i is in the first focal group and $\mathbf{y}_i = (0, 1)$ if examinee i is in the second focal group.

The ability vector of the i th examinee is $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{ir}, \dots, \theta_{iq})^T$ ($i=1, \dots, N$; $r=1, 2, \dots, q$). $\mathbf{a}_j = (a_{j1}, a_{j2})^T$ is the discrimination parameter and d_j is the boundary parameter.

$$\boldsymbol{\gamma}_j = (\gamma_{j1\cdot}, \gamma_{j2\cdot})^T = \begin{pmatrix} \gamma_{j11} & \gamma_{j12} & \dots & \gamma_{j1r} & \dots & \gamma_{j1q} \\ \gamma_{j21} & \gamma_{j22} & \dots & \gamma_{j2r} & \dots & \gamma_{j2q} \end{pmatrix} (r = 1, \dots, q)$$

is the non-uniform DIF parameter, where $\boldsymbol{\gamma}_{j1\cdot}$ is the non-uniform DIF parameter for the first focal group and $\boldsymbol{\gamma}_{j2\cdot}$ is the non-uniform DIF parameter for the second focal group. $\boldsymbol{\beta}_j = (\beta_{j1}, \beta_{j2})^T$ is the uniform DIF parameter, where β_{j1} is the uniform DIF parameter for the first focal group and β_{j2} is the uniform DIF parameter for the second focal group. If item j does not have DIF, then $\boldsymbol{\gamma}_j = \mathbf{0}$ and $\boldsymbol{\beta}_j = \mathbf{0}$. If item j has uniform DIF, then $\boldsymbol{\gamma}_j = \mathbf{0}$.

Suppose the prior distribution of $\boldsymbol{\theta}_i$ in group y is multivariate normal distribution with mean vector of $\boldsymbol{\mu}_y$ and covariance matrix of $\boldsymbol{\Sigma}_y$. The prior density of $\boldsymbol{\theta}_i$ is

$$f(\boldsymbol{\theta}_i | \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}_y|^{-\frac{1}{2}} e^{-0.5(\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)^T |\boldsymbol{\Sigma}_y|^{-1} (\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)}.$$

If i is in $1, \dots, N_1$, then $y_i = (0, 0)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_1$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_1$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$.

If i is in $N_1 + 1, \dots, N_1 + N_2$, then $y_i = (1, 0)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_2$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_2$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$.

If i is in $N_1 + N_2 + 1, \dots, N_1 + N_2 + N_3$, then $y_i = (0, 1)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_3$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_3$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$.

We set the reference group to have zero means and unit variances, that is, $\boldsymbol{\mu}_1 = \mathbf{0}$ and $\text{diag}(\boldsymbol{\Sigma}_1) = \mathbf{1}$. Then with some anchor items, the trait parameters for focal groups, i.e. $\boldsymbol{\mu}_2, \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_2$ and $\boldsymbol{\Sigma}_3$, can be freely estimated.

3 Graded Response Model with DIF

Assume a total sample size N , test length m , number of response categories p and trait dimension q . For a polytomously scored item j , the probability that examinee i with ability vector $\boldsymbol{\theta}_i$ reaching level k or higher on item j is

$$P_{ijk}^* = \frac{1}{1 + e^{-(\mathbf{a}_j \boldsymbol{\theta}_i + d_{jk} + (\mathbf{y}_i \boldsymbol{\gamma}_j) \boldsymbol{\theta}_i + \mathbf{y}_i \boldsymbol{\beta}_{jk})}} (i = 1, \dots, N; j = 1, 2, \dots, m; k = 1, 2, \dots, p-1). \quad (2)$$

\mathbf{y}_i is a group indicator including all the grouping information related to DIF. $\mathbf{y}_i = (0, 0)$ if examinee i is in the reference group, $\mathbf{y}_i = (1, 0)$ if examinee i is in the first focal group and $\mathbf{y}_i = (0, 1)$ if examinee i is in the second focal group.

$$P_{ijk} = P_{ij,k-1}^* - P_{ijk}^* \quad (3)$$

is the probability of an examinee i with ability vector $\boldsymbol{\theta}_i$ reaching response level k on item j .

The trait variable has q dimensions. The ability vector of the i th examinee is $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{ir}, \dots, \theta_{iq})^T$ ($i=1, \dots, N$; $r=1, 2, \dots, q$), and the item parameter matrices are

discrimination parameter

$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_j, \dots, \mathbf{a}_m)^T = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1r} & \dots & a_{1q} \\ a_{21} & a_{22} & \dots & a_{2r} & \dots & a_{2q} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{j1} & a_{j2} & \dots & a_{jr} & \dots & a_{jq} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mr} & \dots & a_{mq} \end{pmatrix} (j = 1, 2, \dots, m; r = 1, \dots, q),$$

boundary parameter

$$\mathbf{D} = (\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_j, \dots, \mathbf{d}_m)^T = \begin{pmatrix} d_{11} & d_{12} & \dots & d_{1k} & \dots & d_{1,p-1} \\ d_{21} & d_{22} & \dots & d_{2k} & \dots & d_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ d_{j1} & d_{j2} & \dots & d_{jk} & \dots & d_{j,p-1} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ d_{m1} & d_{m2} & \dots & d_{mk} & \dots & d_{m,p-1} \end{pmatrix} (j = 1, 2, \dots, m; k = 1, \dots, p-1),$$

non-uniform DIF parameter

$$\mathbf{\Gamma} = (\gamma_1, \gamma_2, \dots, \gamma_j, \dots, \gamma_m) (j = 1, \dots, m)$$

$$\gamma_j = (\gamma_{j1\cdot}, \gamma_{j2\cdot})^T = \begin{pmatrix} \gamma_{j11} & \gamma_{j12} & \dots & \gamma_{j1r} & \dots & \gamma_{j1q} \\ \gamma_{j21} & \gamma_{j22} & \dots & \gamma_{j2r} & \dots & \gamma_{j2q} \end{pmatrix} (r = 1, \dots, q),$$

where q is the dimension of $\boldsymbol{\theta}$, and each row of γ_j is (non-uniform) DIF parameter for a focal group, i.e. $\gamma_{j1\cdot}$ is (non-uniform) DIF parameter for the first focal group, and $\gamma_{j2\cdot}$ is (non-uniform) DIF parameter for the second focal group,

and uniform DIF parameter

$$\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_j, \dots, \beta_m) (j = 1, \dots, m)$$

$$\beta_j = (\beta_{j1\cdot}, \beta_{j2\cdot})^T = \begin{pmatrix} \beta_{j11} & \beta_{j12} & \dots & \beta_{j1k} & \dots & \beta_{j1,p-1} \\ \beta_{j21} & \beta_{j22} & \dots & \beta_{j2k} & \dots & \beta_{j2,p-1} \end{pmatrix} (k = 1, \dots, p-1),$$

where p is the number of categories in GRM, and each row of β_j is (uniform) DIF parameter for a focal group, i.e. $\beta_{j1\cdot}$ is (uniform) DIF parameter for the first focal group, and $\beta_{j2\cdot}$ is (uniform) DIF parameter for the second focal group.

If an item does not have DIF, then $\mathbf{\Gamma} = \mathbf{0}$ and $\boldsymbol{\beta} = \mathbf{0}$. If an item has uniform DIF, then $\mathbf{\Gamma} = \mathbf{0}$.

The $N * m$ response matrix is

$$\mathbf{U} = \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1j} & \dots & u_{1m} \\ u_{21} & u_{22} & \dots & u_{2j} & \dots & u_{2m} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ u_{i1} & u_{i2} & \dots & u_{ij} & \dots & u_{im} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ u_{N1} & u_{N2} & \dots & u_{Nj} & \dots & u_{Nm} \end{pmatrix} (i = 1, \dots, N; j = 1, 2, \dots, m).$$

A dummy variable to indicate whether examinee i gets score k for the item j

$$x_{ijk} = \begin{cases} 1, & \text{if } u_{ij} = k \\ 0, & \text{if } u_{ij} \neq k \end{cases}.$$

\mathbf{y}_i is the group indicator. $\mathbf{y}_i = (0, 0)$ stands for the reference group, $\mathbf{y}_i = (1, 0)$ stands for the first focal group and $\mathbf{y}_i = (0, 1)$ stands for the second focal group. The sample size of the reference group, the first focal group and the second focal group are denoted by N_1 , N_2 , N_3 , respectively. We have the total sample size $N = N_1 + N_2 + N_3$. We have

$$\mathbf{Y} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_{N_1} \\ \mathbf{y}_{N_1+1} \\ \mathbf{y}_{N_1+2} \\ \vdots \\ \mathbf{y}_{N_1+N_2} \\ \mathbf{y}_{N_1+N_2+1} \\ \mathbf{y}_{N_1+N_2+2} \\ \vdots \\ \mathbf{y}_{N_1+N_2+N_3} \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & 1 \end{pmatrix}.$$

Suppose the prior distribution of $\boldsymbol{\theta}_i$ in group y is multivariate normal distribution with mean vector of $\boldsymbol{\mu}_y$ and covariance matrix of $\boldsymbol{\Sigma}_y$. The prior density of $\boldsymbol{\theta}_i$ is

$$f(\boldsymbol{\theta}_i | \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}_y|^{-\frac{1}{2}} e^{-0.5(\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)^T |\boldsymbol{\Sigma}_y|^{-1} (\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)}.$$

If i is in $1, \dots, N_1$, then $y_i = (0, 0)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_1$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_1$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$.

If i is in $N_1 + 1, \dots, N_1 + N_2$, then $y_i = (1, 0)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_2$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_2$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$.

If i is in $N_1 + N_2 + 1, \dots, N_1 + N_2 + N_3$, then $y_i = (0, 1)$. $\boldsymbol{\mu}_y = \boldsymbol{\mu}_3$ and $\boldsymbol{\Sigma}_y = \boldsymbol{\Sigma}_3$, $\boldsymbol{\theta}_i \sim N(\boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)$.

We have $\boldsymbol{\mu}_1 = \mathbf{0}$ and all elements on the diagonal of $\boldsymbol{\Sigma}_1$ are 1 for the reference group. Then with some anchor items, the trait parameters for focal groups, i.e. $\boldsymbol{\mu}_2$, $\boldsymbol{\mu}_3$, $\boldsymbol{\Sigma}_2$ and $\boldsymbol{\Sigma}_3$, can be freely estimated.

Denote G_0 as the number of points we evenly take from each coordinate dimension. Then $G = G_0^q$ quadrature samples (same for all examinees) are denoted by $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_g, \dots, \mathbf{q}_G)^T$ ($g = 1, \dots, G$), and $\mathbf{q}_g = (q_{g1}, q_{g2}, \dots, q_{gr}, \dots, q_{gq})$ ($r=1, 2, \dots, q$). At iteration t , we calculate $f(\mathbf{q}_g | \boldsymbol{\mu}_y^{(t-1)}, \boldsymbol{\Sigma}_y^{(t-1)})$ for each group y , where $\boldsymbol{\mu}_y^{(t-1)}$ and $\boldsymbol{\Sigma}_y^{(t-1)}$ are the estimated trait parameters from last iteration.

For an examinee i in the reference group (group 1), we have $y = 1$ and

$$P_{ijk|q_g}^* = P_{jyk|q_g}^* = P_{j1k|q_g}^* = \frac{1}{1 + e^{-(\mathbf{a}_j \mathbf{q}_g + d_k)}} \\ (i = 1, \dots, N_1; j = 1, 2, \dots, m; k = 1, 2, \dots, p-1; g = 1, \dots, G).$$

For an examinee i in the first focal group (group 2), $y = 2$ and

$$P_{ijk|q_g}^* = P_{jyk|q_g}^* = P_{j2k|q_g}^* = \frac{1}{1 + e^{-(\mathbf{a}_j \mathbf{q}_g + d_k + \gamma_{j1} \cdot \mathbf{q}_g + \beta_{j1k})}}$$

$$(i = N_1 + 1, \dots, N_1 + N_2; j = 1, 2, \dots, m; k = 1, 2, \dots, p - 1; g = 1, \dots, G).$$

For the second focal group (group 3), $y = 3$ and

$$P_{ijk|q_g}^* = P_{jyk|q_g}^* = P_{j3k|q_g}^* = \frac{1}{1 + e^{-(\mathbf{a}_j \mathbf{q}_g + d_k + \gamma_{j2} \cdot \mathbf{q}_g + \beta_{j2k})}}.$$

$$(i = N_1 + N_2 + 1, \dots, N_1 + N_2 + N_3; j = 1, 2, \dots, m; k = 1, 2, \dots, p - 1; g = 1, \dots, G).$$

$$P_{jyk|q_g} = P_{j,y,k-1|q_g}^* - P_{jyk|q_g}^*$$

4 Model Identifiability Constraint

Yet the model is not identified. Some constraints on the item parameters are required. Here, for each dimension, we set one anchor item which we know its DIF parameters (Γ and β) are zero for all groups.

For instance, if we have two ability dimensions ($q=2$), test length $m = 20$, and the each factor is loaded on 10 items, then the simple structure discrimination parameter matrix will take the form

$$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m)^T = \begin{pmatrix} a_{11} & 0 \\ 0 & a_{22} \\ a_{31} & 0 \\ a_{41} & 0 \\ a_{51} & 0 \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ a_{10,1} & 0 \\ a_{11,1} & 0 \\ 0 & a_{12,2} \\ 0 & a_{13,2} \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \\ 0 & a_{19,2} \\ 0 & a_{20,2} \end{pmatrix},$$

and the DIF parameters are

$$\mathbf{\Gamma} = (\mathbf{0}, \mathbf{0}, \mathbf{\Gamma}_3, \dots, \mathbf{\Gamma}_m)$$

and

$$\mathbf{\beta} = (\mathbf{0}, \mathbf{0}, \mathbf{\beta}_3, \dots, \mathbf{\beta}_m).$$

We further assume the reference group has mean zero and variance one and only estimate its correlation, and the means and all the elements in covariance matrices of two focal groups can be freely estimated.

5 Uniform DIF Detection via LASSO

As mentioned before, if an item has uniform DIF, then $\boldsymbol{\Gamma} = \mathbf{0}$. The DIF parameter we are estimating is only $\boldsymbol{\beta} = (\mathbf{0}, \dots, \mathbf{0}, \boldsymbol{\beta}_{q+1}, \dots, \boldsymbol{\beta}_m)$.

5.1 E step

For an examinee with ability $\boldsymbol{\theta}_i$ the conditional likelihood of observing \mathbf{u}_i is

$$L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\theta}_i \mid \mathbf{y}, \mathbf{u}_i) = \prod_{j=1}^m \prod_{k=1}^p P_{jk}(\boldsymbol{\theta}_i)^{x_{ijk}}. \quad (4)$$

With the assumption of prior distribution of latent trait, the joint likelihood of \mathbf{u}_i and $\boldsymbol{\theta}_i$ is

$$\begin{aligned} L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{y}, \mathbf{u}_i, \boldsymbol{\theta}_i) &= L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\theta}_i \mid \mathbf{y}, \mathbf{u}_i) f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \boldsymbol{\theta}_i) \\ &= \prod_{j=1}^m \prod_{k=1}^p P_{jk}(\boldsymbol{\theta}_i)^{x_{ijk}} (2\pi)^{-p/2} |\boldsymbol{\Sigma}_y|^{-1/2} \exp(-0.5(\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)' \boldsymbol{\Sigma}_y^{-1} (\boldsymbol{\theta}_i - \boldsymbol{\mu}_y)). \end{aligned} \quad (5)$$

Therefore, the marginal likelihood of \mathbf{u}_i is

$$m(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{y}, \mathbf{u}_i) = \int L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta} \mid \mathbf{y}, \mathbf{u}_i, \boldsymbol{\theta}_i) f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \boldsymbol{\theta}_i) d\boldsymbol{\theta}_i \quad (6)$$

Then

$$h(\boldsymbol{\theta}_i \mid \mathbf{u}_i, \mathbf{y}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_y^{(t-1)}, \boldsymbol{\Sigma}_y^{(t-1)}) = \frac{L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta} \mid \mathbf{y}, \mathbf{u}_i, \boldsymbol{\theta}_i) f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \boldsymbol{\theta}_i)}{m(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{y}, \mathbf{u}_i)} \quad (7)$$

is the posterior density of $\boldsymbol{\theta}_i$ given the estimation of \mathbf{A} , \mathbf{D} , $\boldsymbol{\beta}$ and $\boldsymbol{\Sigma}$ at the iteration t .

The expected complete data log-likelihood with respect to the posterior distribution of $\boldsymbol{\theta}$

$$\begin{aligned} &E[\log\{L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{Y}, \mathbf{U}, \boldsymbol{\Theta})\} \mid \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \mathbf{Y}, \mathbf{U}, \boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)}] \\ &= \sum_i^N \left\{ \int \log L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta} \mid \mathbf{y}, \mathbf{u}_i, \boldsymbol{\theta}_i) h(\boldsymbol{\theta}_i \mid \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)}) d\boldsymbol{\theta}_i \right. \\ &\quad \left. + \int \log f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \boldsymbol{\theta}_i) h(\boldsymbol{\theta}_i \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_y^{(t-1)}, \boldsymbol{\Sigma}_y^{(t-1)}) d\boldsymbol{\theta}_i \right\} \end{aligned} \quad (8)$$

At iteration t , applying Gauss-Hermite quadrature nodes and the integration above can be updated as

$$\begin{aligned}
& E[\log L(\mathbf{A}, \mathbf{D}, \beta, \mu, \Sigma \mid \mathbf{Y}, \mathbf{U})] \\
&= \sum_i^N \sum_g^G \log L(\mathbf{A}, \mathbf{D}, \beta \mid \mathbf{u}_i, \mathbf{q}_g) \frac{L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)} \\
&+ \sum_i^N \sum_g^G \log f(\mu, \Sigma \mid \mathbf{q}_g) \frac{L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)} \\
&= \sum_i^N \sum_g^G \sum_j^m \sum_k^p x_{ijk} \log P_{ijk|\mathbf{q}_g} \frac{L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)} \\
&+ \sum_i^N \sum_g^G \log f(\mu, \Sigma \mid \mathbf{q}_g) \frac{L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g \mid \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mathbf{Y}, \mu^{(t-1)}, \Sigma^{(t-1)}) f(\mu^{(t-1)}, \Sigma^{(t-1)} \mid \mathbf{q}_g)}
\end{aligned} \tag{9}$$

Then we can define two artificial terms.

For the reference group, $y = 1$. We have

$$n_{gy} = n_{g1} = \sum_{i=1}^{N_1} \frac{L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_1^{(t-1)}, \boldsymbol{\Sigma}_1^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_1^{(t-1)}, \boldsymbol{\Sigma}_1^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)},$$

and

$$r_{gjjk} = r_{gjl k} = \sum_{i=1}^{N_1} x_{ijk} \frac{L(q_g | y_i, u_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_1^{(t-1)}, \boldsymbol{\Sigma}_1^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | q_g)}{\sum_a^G L(q_a | y_i, u_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_1^{(t-1)}, \boldsymbol{\Sigma}_1^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | q_a)}.$$

For the first focal group, $y = 2$. We have

$$n_{gy} = n_{g2} = \sum_{i=N_1+1}^{N_1+N_2} \frac{L(\mathbf{q}_g \mid \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_2^{(t-1)}, \boldsymbol{\Sigma}_2^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} \mid \mathbf{q}_g)}{\sum_{i=N_1+1}^G L(\mathbf{q}_g \mid \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_2^{(t-1)}, \boldsymbol{\Sigma}_2^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} \mid \mathbf{q}_g)},$$

and

$$r_{gjjk} = r_{gj2k} = \sum_{i=N_1+1}^{N_1+N_2} x_{ijk} \frac{L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_2^{(t-1)}, \boldsymbol{\Sigma}_2^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)}{\sum_g^G L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_2^{(t-1)}, \boldsymbol{\Sigma}_2^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)}.$$

For the second focal group, $y = 3$. We have

$$n_{gy} = n_{g3} = \sum_{i=N_1+N_2+1}^{N_1+N_2+N_3} \frac{L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_3^{(t-1)}, \boldsymbol{\Sigma}_3^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)}{\sum_a^G L(\mathbf{q}_g | \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_3^{(t-1)}, \boldsymbol{\Sigma}_3^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} | \mathbf{q}_g)}$$

and

$$r_{gjjk} = r_{gj3k} = \sum_{i=N_1+N_2+1}^{N_1+N_2+N_3} x_{ijk} \frac{L(\mathbf{q}_g \mid \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_3^{(t-1)}, \boldsymbol{\Sigma}_3^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} \mid \mathbf{q}_g)}{\sum_a^G L(\mathbf{q}_a \mid \mathbf{y}_i, \mathbf{u}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \boldsymbol{\beta}^{(t-1)}, \boldsymbol{\mu}_3^{(t-1)}, \boldsymbol{\Sigma}_3^{(t-1)}) f(\boldsymbol{\mu}^{(t-1)}, \boldsymbol{\Sigma}^{(t-1)} \mid \mathbf{q}_a)}.$$

$n_g = n_{g1} + n_{g2} + n_{g3}$ represents the expected number of examinees with the ability \mathbf{q}_g , and $r_{jgk} = r_{jgk1} + r_{jgk2} + r_{jgk3}$ is the expected number of examinees who get the score level k on the item j with the ability \mathbf{q}_g .

$$E[\log\{L(\mathbf{A}, \mathbf{D}, \boldsymbol{\beta}, \boldsymbol{\mu}, \boldsymbol{\Sigma} \mid \mathbf{Y}, \mathbf{U}, \boldsymbol{\Theta})\}] = \sum_g^G \sum_j^m \sum_y^3 \sum_k^p (r_{gjk} \log P_{jyk|q_g}) + \sum_g^G \sum_y^3 n_{gy} \log f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{q}_g) \quad (10)$$

In the EM problem, we want to maximize the above expectation at the iteration t . Denote this unpenalized expectation as $\log M$.

For each item j , we define

$$\log M_j = \sum_g^G \sum_y^3 \sum_k^p (r_{jgk} \log P_{jky|q_g}) + \sum_g^G \sum_y^3 n_{gy} \log f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{q}_g) \quad (11)$$

In our uniform DIF detection problem, the maximum likelihood method does not serve the purpose of DIF variable selection. We apply lasso and minimize the following objective function

$$-\log M + \eta \sum_j^m \|\boldsymbol{\beta}_j\|_1 \quad (12)$$

For each item, we minimize

$$-\log M_j + \eta \|\boldsymbol{\beta}_j\|_1 \quad (13)$$

where η is the lasso tuning parameter.

$$(\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\boldsymbol{\beta}}) = \operatorname{argmin}\{-\log M + \eta \|\boldsymbol{\beta}\|_1\} \quad (14)$$

5.2 M step

In our DIF detection problem, we assume the reference group has mean zero and variance one and only estimate the correlation, and the means and all the elements in covariance matrices of two focal groups can be freely estimated.

In quadrature method, at the iteration t , the first partial derivative with respect to $\boldsymbol{\mu}$ is

$$\begin{aligned} \frac{\partial \log M}{\partial \boldsymbol{\mu}_y} &= \sum_g^G n_{gy} \frac{\partial \log f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{q}_g)}{\partial \boldsymbol{\mu}_y} \\ &= \sum_g^G n_{gy} \frac{\partial -\frac{1}{2}(\mathbf{q}_g - \boldsymbol{\mu}_y)^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{q}_g - \boldsymbol{\mu}_y)}{\partial \boldsymbol{\mu}_y} \\ &= \sum_g^G n_{gy} (\mathbf{q}_g - \boldsymbol{\mu}_y) \boldsymbol{\Sigma}_y^{-1} \end{aligned} \quad (15)$$

Set $\frac{\partial \log M}{\partial \boldsymbol{\mu}_y} = 0$, and we know that $\sum_g^G n_{gy} = N_y$.

$\hat{\boldsymbol{\mu}}_y$ can be updated as

$$\hat{\boldsymbol{\mu}}_2 = \frac{\sum_{g=1}^G n_{g2} \mathbf{q}_g}{N_2}, \quad (16)$$

and

$$\hat{\boldsymbol{\mu}}_3 = \frac{\sum_{g=1}^G n_{g3} \mathbf{q}_g}{N_3}. \quad (17)$$

The first partial derivative with respect to $\boldsymbol{\Sigma}$ is

$$\begin{aligned} \frac{\partial \log M}{\partial \boldsymbol{\Sigma}_y} &= \sum_g^G n_{gy} \frac{\partial \log f(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y \mid \mathbf{q}_g)}{\partial \boldsymbol{\Sigma}_y} \\ &= \sum_g^G n_{gy} \frac{\partial (-\frac{g}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_y| - \frac{1}{2} (\mathbf{q}_g - \boldsymbol{\mu}_y)^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{q}_g - \boldsymbol{\mu}_y))}{\partial \boldsymbol{\Sigma}_y} \\ &= \sum_g^G n_{gy} [-\frac{1}{2} \boldsymbol{\Sigma}_y^{-1} + \frac{1}{2} \boldsymbol{\Sigma}_y^{-1} (\mathbf{q}_g - \boldsymbol{\mu}_y) (\mathbf{q}_g - \boldsymbol{\mu}_y)^T \boldsymbol{\Sigma}_y^{-1}] \end{aligned} \quad (18)$$

Set $\frac{\partial \log M}{\partial \boldsymbol{\mu}_y} = 0$, and use the fact that $\sum_g^G n_{gy} = N_y$.

$\hat{\boldsymbol{\Sigma}}_y$ can be updated as

$$\hat{\boldsymbol{\Sigma}}_1 = \frac{\sum_{g=1}^G n_{g1} \mathbf{q}_g \mathbf{q}_g'}{N_1}, \quad (19)$$

$$\hat{\boldsymbol{\Sigma}}_2 = \frac{\sum_{g=1}^G n_{g2} (\mathbf{q}_g - \hat{\boldsymbol{\mu}}_2) (\mathbf{q}_g - \hat{\boldsymbol{\mu}}_2)'}{N_2}, \quad (20)$$

and

$$\hat{\boldsymbol{\Sigma}}_3 = \frac{\sum_{g=1}^G n_{g3} (\mathbf{q}_g - \hat{\boldsymbol{\mu}}_3) (\mathbf{q}_g - \hat{\boldsymbol{\mu}}_3)'}{N_3}. \quad (21)$$

To standardize the covariance matrix, we calculate standardized quadrature points for the later steps.

$$\mathbf{q}_g^* = \frac{\mathbf{q}_g}{\sqrt{\text{diag} \hat{\boldsymbol{\Sigma}}_1}}. \quad (22)$$

Then we do the following transformation on mean vector and covariance matrices for three groups.

$$\hat{\boldsymbol{\Sigma}}_1^* = \frac{\sum_{g=1}^G n_{g1} \mathbf{q}_g^* \mathbf{q}_g^{*'}}{N_1}, \quad (23)$$

$$\hat{\boldsymbol{\Sigma}}_2^* = \frac{\sum_{g=1}^G n_{g2} (\mathbf{q}_g^* - \hat{\boldsymbol{\mu}}_2) (\mathbf{q}_g^* - \hat{\boldsymbol{\mu}}_2)'}{N_2}, \quad (24)$$

and

$$\hat{\Sigma}_3^* = \frac{\sum_{g=1}^G n_{g3}(\mathbf{q}_g^* - \hat{\mu}_3)(\mathbf{q}_g^* - \hat{\mu}_3)'}{N_3}. \quad (25)$$

the first partial derivative with respect to a_{jr} is

$$\frac{\partial \log M}{\partial a_{jr}} = \sum_{g=1}^G \sum_y^3 \sum_{k=1}^p \left(\frac{r_{gjjk} q_{gr}}{P_{jky|q_g}} (\omega_{j,y,k-1} - \omega_{jyk}) \right) \quad (26)$$

where $\omega_{jyk} = P_{jyk|q_g}^* - (P_{jky|q_g}^*)^2$.

Similarly, we have the first partial derivative with respect to d_{jk}

$$\frac{\partial \log M}{\partial d_{jk}} = \sum_g^G \sum_y^3 \omega_{jky} \left(\frac{r_{gjj,(k+1),y}}{P_{jy,(k+1)|q_g}} - \frac{r_{gjjk}}{P_{jyk|q_g}} \right) \quad (27)$$

where $\omega_{jyk} = P_{jyk|q_g}^* - (P_{jky|q_g}^*)^2$,

and the first partial derivative with respect to β_{jyk} , where $y=(2,3)$, is

$$\frac{\partial \log M}{\partial \beta_{jyk}} = \sum_g^G \omega_{jyk} \left(\frac{r_{gjj,(k+1)}}{P_{jy,(k+1)|q_g}} - \frac{r_{gjjk}}{P_{jyk|q_g}} \right) \quad (28)$$

where $\omega_{jyk} = P_{jyk|q_g}^* - (P_{jky|q_g}^*)^2$.

The second partial derivatives in the Hessian matrix are given by

$$\begin{aligned} \frac{\partial^2 \log M}{\partial a_{jr}^2} &= \sum_{g=1}^G \sum_y^3 \sum_{k=1}^p - \frac{r_{gjjk} q_{gr}^2 (P_{jy,(k-1)|q_g}^* Q_{jy,(k-1)|q_g}^* - P_{jyk|q_g}^* Q_{jyk|q_g}^*)^2}{P_{jyk|q_g}^2} \\ &= \sum_{g=1}^G \sum_y^3 \sum_{k=1}^p - \frac{r_{gjjk} q_{gr}^2 (\omega_{jy(k-1)} - \omega_{jyk})}{P_{jyk|q_g}^2} \end{aligned} \quad (29)$$

$$\begin{aligned} \frac{\partial^2 \log M}{\partial d_{jk}^2} &= \sum_y^3 \sum_{g=1}^G - \left(\frac{r_{gjjk}}{P_{jyk|q_g}^2} + \frac{r_{gjj(k+1)}}{P_{jy(k+1)|q_g}^2} \right) P_{jyk|q_g}^{*2} (1 - P_{jyk|q_g}^*)^2 \\ &= \sum_y^3 \sum_{g=1}^G - \left(\frac{r_{gjjk}}{P_{jyk|q_g}^2} + \frac{r_{gjj(k+1)}}{P_{jy(k+1)|q_g}^2} \right) \omega_{jyk}^2 \end{aligned} \quad (30)$$

$$\begin{aligned} \frac{\partial^2 \log M}{\partial d_{jk} \partial d_{j,k+1}} &= \sum_{g=1}^G \sum_y^3 \frac{r_{gjj(k+1)}}{P_{jy(k+1)|q_g}^2} (P_{jyk|q_g}^{*2} (1 - P_{jyk|q_g}^*)^2) (P_{jy(k+1)|q_g}^{*2} (1 - P_{jy(k+1)|q_g}^*)^2) \\ &= \sum_{g=1}^G \sum_y^3 \frac{r_{gjj(k+1)}}{P_{jy(k+1)|q_g}^2} \omega_{jyk}^2 \omega_{jy(k+1)}^2 \end{aligned} \quad (31)$$

and

$$\begin{aligned}
\frac{\partial^2 \log M}{\partial a_{jr} \partial d_{jk}} &= \sum_{g=1}^G \sum_y^3 P_{jyk}^* Q_{jyk}^* q_{gr} \left[\frac{r_{gjjyk}}{P_{jyk|q_g}^2} (P_{jy(k-1)|q_g}^* Q_{jy(k-1)|q_g}^* - P_{jyk|q_g}^* Q_{jyk|q_g}^*) \right. \\
&\quad \left. + \frac{r_{gjjy(k+1)}}{P_{jy(k+1)|q_g}^2} (P_{jyk|q_g}^* Q_{jyk|q_g}^* - P_{jy(k+1)|q_g}^* Q_{jy(k+1)|q_g}^*) \right] \\
&= \sum_{g=1}^G \sum_y^3 \omega_{jyk} q_{gr} \left[\frac{r_{gjjyk}}{P_{jyk|q_g}^2} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{r_{gjjy(k+1)}}{P_{jy(k+1)|q_g}^2} (\omega_{jyk} - \omega_{jy(k+1)}) \right]
\end{aligned} \tag{32}$$

where

$$\begin{aligned}
Q_{jyk|q_g}^* &= 1 - P_{jyk|q_g}^* \\
\omega_{jyk} &= P_{jyk|q_g}^* * Q_{jyk|q_g}^*
\end{aligned}$$

$$\frac{\partial^2 \log M}{\partial \beta_{jyk}^2} = \frac{\partial^2 \log M}{\partial \beta_{jyk} \partial d_{jk}} = \sum_{g=1}^G - \left(\frac{r_{gjjyk}}{P_{jyk|q_g}^2} + \frac{r_{gjjy(k+1)}}{P_{jy(k+1)|q_g}^2} \right) P_{jyk|q_g}^{*2} (1 - P_{jyk|q_g}^*)^2 \tag{33}$$

$$\frac{\partial^2 \log M}{\partial a_{jr} \partial \beta_{jyk}} = \sum_{g=1}^G \omega_{jyk} q_{gr} \left[\frac{r_{gjjyk}}{P_{jyk|q_g}^2} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{r_{gjjy(k+1)}}{P_{jy(k+1)|q_g}^2} (\omega_{jyk} - \omega_{jy(k+1)}) \right] \tag{34}$$

The expectation of the second partial derivatives in the Fisher scoring method are given by

$$E\left(\frac{\partial^2 \log M}{\partial a_{jr}^2}\right) = \sum_{g=1}^G \sum_y^3 \sum_{k=1}^p - \frac{n_{gy} q_{gr}^2 (\omega_{jy(k-1)} - \omega_{jyk})}{P_{jyk|q_g}}, \tag{35}$$

$$E\left(\frac{\partial^2 \log M}{\partial d_{jk}^2}\right) = \sum_{g=1}^G \sum_y^3 -n_{gy} \left(\frac{1}{P_{jyk|q_g}} + \frac{1}{P_{jy(k+1)|q_g}} \right) \omega_{jyk}^2, \tag{36}$$

$$E\left(\frac{\partial^2 \log M}{\partial d_{jk} \partial d_{j,k+1}}\right) = \sum_{g=1}^G \sum_y^3 \frac{n_{gy}}{P_{jy(k+1)|q_g}} \omega_{jyk}^2 \omega_{jy(k+1)}^2, \tag{37}$$

and

$$E\left(\frac{\partial^2 \log M}{\partial a_{jr} \partial d_{jk}}\right) = \sum_{g=1}^G \sum_y^3 n_{gy} \omega_{jyk} q_{gr} \left[\frac{1}{P_{jyk|q_g}} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{1}{P_{jy(k+1)|q_g}} (\omega_{jyk} - \omega_{jy(k+1)}) \right]. \tag{38}$$

$$E\left(\frac{\partial^2 \log M}{\partial \beta_{jyk}^2}\right) = E\left(\frac{\partial^2 \log M}{\partial \beta_{jyk} \partial d_{jk}}\right) = \sum_{g=1}^G -n_{gy} \left(\frac{1}{P_{jyk|q_g}} + \frac{1}{P_{jy(k+1)|q_g}} \right) \omega_{jyk}^2. \tag{39}$$

$$E\left(\frac{\partial^2 \log M}{\partial a_{jr} \partial \beta_{jyk}}\right) = \sum_{g=1}^G n_{gy} \omega_{jyk} q_{gr} \left[\frac{1}{P_{jyk|q_g}} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{1}{P_{jy(k+1)|q_g}} (\omega_{jyk} - \omega_{jy(k+1)}) \right]. \tag{40}$$

5.2.1 Cyclical coordinate descent

By Bazaraa, Sherali, and Shetty (2006), for a convex function f , a point $\bar{\theta}$ is a global minimizer of f if and only if $\partial f(\bar{\theta})$, the subgradient of f at $\bar{\theta}$, contains 0. Hence $\hat{\theta}_\tau$ is the global minimizer only when $\hat{\theta}_\tau = \text{sign}(s)(|s| - \tau)_+$, where $(u)_+ = u1(u > 0)$. This is called the soft-threshold of s and τ , and can be denoted by

$$\begin{aligned}\hat{\theta}_\tau &= \text{soft}(s, \tau) \equiv \text{sign}(s)(|s| - \tau)_+ \\ &= \arg \min_{\theta \in \mathbb{R}} \{0.5\theta^2 - s\theta + \tau|\theta|\}.\end{aligned}\tag{41}$$

Then, to minimize our objective function with respect to β , we calculate a second-order Tylor approximation of the marginal log-likelihood at $\beta^{(t-1)}$, and our lasso estimator in (13) can be updated by

$$\begin{aligned}\hat{\beta} &= \text{argmin}\{-\log M(\beta) + \eta\|\beta\|_1\} \\ &= \text{argmin}\{-\log M(\beta^{(t-1)}) - \partial_\beta \log M(\beta^{(t-1)})(\beta - \beta^{(t-1)}) - \frac{\partial_\beta^2 \log M(\beta^{(t-1)})}{2}(\beta - \beta^{(t-1)})^2 + \eta\|\beta\|_1\} \\ &= -\frac{\text{soft}(\partial_\beta \log M - \beta_j^{(t-1)} * \partial_\beta^2 \log M, \eta)}{\partial_\beta^2 \log M}\end{aligned}\tag{42}$$

We run an EM and cyclical coordinate descent algorithm given by following.

Algorithm 1: Uniform DIF Detection via LASSO

Input : $A_0, D_0, \beta_0, \mu_0, \Sigma_0, U, \eta, \varepsilon_1, \varepsilon_2$

Output: $\hat{A}, \hat{D}, \hat{\beta}, \hat{\mu}, \hat{\Sigma}$

```

1 set  $t_1 = 1$ ,  $\delta^{(t_1-1)} = \text{any value greater than } \varepsilon_1$ 
2 while  $\delta_1^{(t_1-1)} > \varepsilon_1$  do
3   Calculate  $n_{gy}$  and  $r_{gyk}$ 
4   Update  $\mu^{(t_1)}$  and  $\Sigma^{(t_1)}$ 
5   for  $j=1, \dots, m$  do
6     set  $t_2 = 1$ ,  $\delta^{(t_2-1)} = \text{any value greater than } \varepsilon_2$ 
7     while  $\delta_2^{(t_2-1)} > \varepsilon_2$  do
8       Calculate  $P_{jyk|q_g}^*, Q_{jyk|q_g}^*$ 
9        $a_{jr}^{(t_2)} = a_{jr}^{(t_2-1)} - \frac{\partial_{a_{jr}} \log M}{\partial_{a_{jr}}^2 \log M}$ 
10       $d_{jk}^{(t_2)} = d_{jk}^{(t_2-1)} - \frac{\partial_{d_{jk}} \log M}{\partial_{d_{jk}}^2 \log M}$ 
11       $\beta_{jyk}^{(t_2)} = -\frac{\text{soft}(\partial_{\beta_{jyk}} \log M - \beta_{jyk}^{(t_2-1)} * \partial_{\beta_{jyk}}^2 \log M, \eta)}{\partial_{\beta_{jyk}}^2 \log M}$ 
12       $\delta_2^{(t_2)} = \|A_j^{(t_2)} - A_j^{(t_2-1)}\| + \|D_j^{(t_2)} - D_j^{(t_2-1)}\| + \|\beta_j^{(t_2)} - \beta_j^{(t_2-1)}\|$ 
13       $t_2 = t_2 + 1$ 
14    end
15     $a_{jr}^{(t_1)*} = a_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})}$ 
16  end
17   $\delta_1^{(t_1)} = \|A^{(t_1)} - A^{(t_1-1)}\| + \|D^{(t_1)} - D^{(t_1-1)}\| + \|\beta^{(t_1)} - \beta^{(t_1-1)}\|$ 
18   $t_1 = t_1 + 1$ 
19 end
```

Note that $\text{diag}(\hat{\Sigma}_{1r})$ is the r th element on the diagonal of the estimated covariance matrix of the reference group $\hat{\Sigma}_1$.

$\varepsilon_1 = 10^{-3}$ and $\varepsilon_2 = 10^{-7}$.

5.3 Simulation

Sample Size. The total sample size is $N = 1500$, and the group sample sizes are $N_1 = N_2 = N_3 = 500$.

Test Length. $m = 20$. Simple structure. 10 items per dimension.

Proportion of DIF. 4 items with DIF. 2 DIF items per dimension.

Magnitude of DIF. The first focal group with larger difficulty parameters (+0.5) on the 4 DIF items.

The second focal group with much larger difficulty parameters (+1) on the 4 DIF items.

Generated parameters.

$a_{j1} \sim U(1.5, 2.5), j = 1, \dots, 10$

$a_{j2} \sim U(1.5, 2.5), j = 11, \dots, 20$

$d_1 \sim N(0, 1)$

$$\mathbf{A} = \begin{pmatrix} 2.17 & 0 \\ 0 & 2.46 \\ 2.41 & 0 \\ 2.45 & 0 \\ 2.34 & 0 \\ 1.84 & 0 \\ 1.85 & 0 \\ 1.92 & 0 \\ 1.94 & 0 \\ 1.90 & 0 \\ 1.92 & 0 \\ 0 & 2.43 \\ 0 & 1.82 \\ 0 & 2.22 \\ 0 & 1.93 \\ 0 & 1.88 \\ 0 & 1.84 \\ 0 & 2.12 \\ 0 & 2.42 \\ 0 & 2.15 \end{pmatrix},$$

$$D = \begin{pmatrix} 0.03 \\ -1.28 \\ 0.58 \\ -2.06 \\ 0.12 \\ 3.25 \\ -0.41 \\ -0.51 \\ 0.89 \\ 1.33 \\ 0.85 \\ 0.82 \\ -0.37 \\ -0.99 \\ -0.27 \\ 0.19 \\ 1.73 \\ 0.05 \\ -1.86 \\ -0.63 \end{pmatrix},$$

$$\beta = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.5 & 1 \\ 0.5 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.5 & 1 \\ 0.5 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix},$$

No impact. $\theta_i \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.85 \\ 0.85 & 1 \end{pmatrix}\right)$

Tuning parameters. Tuning parameters η are chosen from a sequence starting from 21 to 51 with increment 3.

5.3.1 Results of 50 Replications

Table 1. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.985	0.34	0.985
Type I	0.047	0.024	0.024

Omnibus DIF is defined as if at least one focal group showd DIF on an item, then that item is flagged as DIF.

Table 2. Type I error and Power of mirt LRT (0.05 significance level)

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.86	0.23	0.91
Type I	0.00714	0.0085	0.002857

mirt LRT can only do the omnibus DIF test.

mirt wald

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.2031	0.2031	0.9531
Type I	0.0045	0.0045	0.0134

Table 3. Item parameter estimates by regularization

Item Parameters	a_1	a_2	d
Bias	-0.011145	-0.00684	-0.04356
RMSE	0.1676937	0.1614306	0.1576307

Table 4. Item parameter estimates by mirt LRT (0.05 significance level)

Item Parameters	a_1	a_2	d
Bias	0.01262	0.00636	-0.00532
RMSE	0.17259	0.16686	0.15026

Our regularization method has slightly better non-DIF item parameter estimates.

Table 5. Absolute bias for DIF magnitude recoveries that were true DIF

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Regularization (include false negative)	0.1529	0.3779	0.2284
mirt LRT (include false negative)	0.2222	0.4181	0.2181

Table 6. Absolute bias for DIF magnitude recoveries that were non – DIF

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Regularization	0.336	0.33	0.341
mirt LRT	0.15696	0.1516	0.1622

The results in Table 6 are the average of estimated DIF for false positive items. LRT by mirt performs better when type I error happens. The type I error is low, so the probability to have these bias is low.

6 Non-uniform DIF Detection via LASSO

When the items have non-uniform DIF on slope only, i.e., there is no DIF on the intercepts, the DIF parameter we are estimating is $\mathbf{\Gamma} = (\mathbf{0}, \dots, \mathbf{0}, \mathbf{\Gamma}_{q+1}, \dots, \mathbf{\Gamma}_m)$.

6.1 E step

For the expectation step, we can use the result in Section 3.1 only replacing the DIF parameter β by $\mathbf{\Gamma}$, and minimize the following objective function

$$-\log M + \eta \sum_j^m \|\mathbf{\Gamma}_j\|_1 \quad (43)$$

For each item, we minimize

$$-\log M_j + \eta \|\mathbf{\Gamma}_j\|_1 \quad (44)$$

where η is the lasso tuning parameter.

$$(\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\mathbf{\Gamma}}) = \operatorname{argmin}\{-\log M + \eta \|\mathbf{\Gamma}\|_1\} \quad (45)$$

6.2 M step

Again, we assume the reference group has mean zero and variance one and only estimate its correlations. The means and all elements in the covariance matrices of two focal groups can be freely estimated.

$\hat{\mu}_2, \hat{\mu}_3, \hat{\Sigma}_1^*, \hat{\Sigma}_2^*, \hat{\Sigma}_3^*$ are same in (15), (16), (22), (23) and (24).

The first partial derivative with respect to a_{jr} and d_{jk} are same as in (25) and (26).

the first partial derivative with respect to γ_{jry} , where $y=(2,3)$, is

$$\begin{aligned} \frac{\partial \log M}{\partial \gamma_{jyr}} &= \sum_g^G \sum_k^p \frac{r_{gjjyk} q_{gr} [P_{jy(k-1)|q_g}^* (1 - P_{jy(k-1)|q_g}^*) - P_{jyk|q_g}^* (1 - P_{jyk|q_g}^*)]}{P_{jyk|q_g}} \\ &= \sum_g^G \sum_k^p \left(\frac{r_{gjjyk} q_{gr}}{P_{jyk|q_g}} (\omega_{jy(k-1)} - \omega_{jyk}) \right) \end{aligned} \quad (46)$$

where $\omega_{jyk} = P_{jky|q_g}^* - (P_{jyk|q_g}^*)^2$.

The second partial derivatives in the Hessian matrix $\frac{\partial^2 \log M}{\partial a_{jr}^2}$, $\frac{\partial^2 \log M}{\partial d_{jk}^2}$, $\frac{\partial^2 \log M}{\partial d_{jk} \partial d_{j,k+1}}$ and $\frac{\partial^2 \log M}{\partial a_{jr} \partial d_{jk}}$ are given by (28)-(31) and their expectations are (34)-(37).

$$\frac{\partial^2 \log M}{\partial \gamma_{jyr}^2} = \frac{\partial^2 \log M}{\partial \gamma_{jyr} \partial a_{jr}} = \sum_{g=1}^G \sum_{k=1}^p -\frac{r_{gjjyk} q_{gr}^2 (\omega_{jy(k-1)} - \omega_{jyk})}{P_{jyk|q_g}^2} \quad (47)$$

$$\frac{\partial^2 \log M}{\partial \gamma_{jyr} \partial d_{jk}} = \sum_{g=1}^G \omega_{jky} q_{gr} \left[\frac{r_{gjjk}}{P_{jyk|q_g}^2} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{r_{gjj(k+1)}}{P_{jy(k+1)|q_g}^2} (\omega_{jyk} - \omega_{jy(k+1)}) \right] \quad (48)$$

The expectation of the second partial derivatives in the Fisher scoring method are given by

$$E\left(\frac{\partial^2 \log M}{\partial \gamma_{jyr}^2}\right) = E\left(\frac{\partial^2 \log M}{\partial \gamma_{jyr} \partial a_{jr}}\right) = \sum_{g=1}^G \sum_{k=1}^p -\frac{n_{gy} q_{gr}^2 (\omega_{jy(k-1)} - \omega_{jyk})}{P_{jyk|q_g}}. \quad (49)$$

$$E\left(\frac{\partial^2 \log M}{\partial \gamma_{jyr} \partial d_{jk}}\right) = \sum_{g=1}^G n_{gy} \omega_{jyk} q_{gr} \left[\frac{1}{P_{jyk|q_g}} (\omega_{jy(k-1)} - \omega_{jyk}) + \frac{1}{P_{jy(k+1)|q_g}} (\omega_{jyk} - \omega_{jy(k+1)}) \right]. \quad (50)$$

6.2.1 Cyclical coordinate descent

Same as in 3.2, to minimize our objective function with respect to $\mathbf{\Gamma}$, our lasso estimator in (13) can be written by

$$\begin{aligned} \hat{\mathbf{\Gamma}} &= \operatorname{argmin}\{-\log M(\mathbf{\Gamma}) + \eta \|\mathbf{\Gamma}\|_1\} \\ &= \operatorname{argmin}\{-\log M(\mathbf{\Gamma}_0) - \partial_{\mathbf{\Gamma}} \log M(\mathbf{\Gamma}_0)(\mathbf{\Gamma} - \mathbf{\Gamma}_0) - \frac{\partial_{\mathbf{\Gamma}}^2 \log M(\mathbf{\Gamma}_0)}{2} (\mathbf{\Gamma} - \mathbf{\Gamma}_0)^2 + \eta \|\mathbf{\Gamma}\|_1\} \\ &= -\frac{\operatorname{soft}(\partial_{\mathbf{\Gamma}} \log M - \mathbf{\Gamma}_j^{(t-1)} * \partial_{\mathbf{\Gamma}}^2 \log M, \eta)}{\partial_{\mathbf{\Gamma}}^2 \log M} \end{aligned} \quad (51)$$

We run a cyclical coordinate descent algorithm for each group (item) with all other groups fixed. For item j , our algorithm is given by following.

1. Calculate $P_{jky|q_g}^*$ and $Q_{jky|q_g}^*$.
2. The parameter a_{jr} and d_{jk} can be updated by

$$a_{jr}^{(t)} = a_{jr}^{(t-1)} - \frac{\partial_{a_{jr}} \log M}{\partial_{a_{jr}}^2 \log M},$$

$$d_{jk}^{(t)} = d_{jk}^{(t-1)} - \frac{\partial_{d_{jk}} \log M}{\partial_{d_{jk}}^2 \log M}$$

and

$$\hat{\mathbf{\Gamma}}_{jyr} = -\frac{\operatorname{soft}(\partial_{\mathbf{\Gamma}_{jyr}} \log M - \mathbf{\Gamma}_{jyr}^{(t-1)} * \partial_{\mathbf{\Gamma}_{jyr}}^2 \log M, \eta)}{\partial_{\mathbf{\Gamma}_{jyr}}^2 \log M}$$

Then we update $P_{jyk|q_g}^*$ and $Q_{jyk|q_g}^*$ by plugging in $\hat{\mathbf{A}}, \hat{\mathbf{D}}$ and $\hat{\beta}$ from last coordinate descent cycle and repeat above steps until a convergence criterion is met.

After we get optimizers for item j , we do transformations on all estimates as following

$$a_{jr}^{(t)*} = a_{jr}^{(t)} * \sqrt{\operatorname{diag}(\hat{\Sigma}_{1r})},$$

$$\gamma_{jr}^{(t)*} = \gamma_{jr}^{(t)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})},$$

We run an EM and cyclical coordinate descent algorithm given by following.

Algorithm 2: Non-uniform DIF Detection via LASSO

Input : $\mathbf{A}_0, \mathbf{D}_0, \mathbf{\Gamma}_0, \boldsymbol{\mu}_0, \Sigma_0, \mathbf{U}, \eta, \varepsilon_1, \varepsilon_2$

Output: $\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\mathbf{\Gamma}}, \hat{\boldsymbol{\mu}}, \hat{\Sigma}$

```

1 set  $t_1 = 1$ ,  $\delta_1^{(t_1-1)}$  = any value greater than  $\varepsilon_1$ 
2 while  $\delta_1^{(t_1-1)} > \varepsilon_1$  do
3   Calculate  $n_{gy}$  and  $r_{gjjk}$ 
4   Update  $\boldsymbol{\mu}^{(t_1)}$  and  $\Sigma^{(t_1)}$ 
5   for  $j=1, \dots, m$  do
6     set  $t_2 = 1$ ,  $\delta_2^{(t_2-1)}$  = any value greater than  $\varepsilon_2$ 
7     while  $\delta_2^{(t_2-1)} > \varepsilon_2$  do
8       Calculate  $P_{jyk|q_g}^*, Q_{jyk|q_g}^*$ 
9        $a_{jr}^{(t_2)} = a_{jr}^{(t_2-1)} - \frac{\partial a_{jr} \log M}{\partial^2 a_{jr} \log M}$ 
10       $d_{jk}^{(t_2)} = d_{jk}^{(t_2-1)} - \frac{\partial d_{jk} \log M}{\partial^2 d_{jk} \log M}$ 
11       $\Gamma_{jyr}^{(t_2)} = -\frac{\text{soft}(\partial \Gamma_{jyr} \log M - \Gamma_{jyr}^{(t_2-1)} * \partial^2 \Gamma_{jyr} \log M, \eta)}{\partial^2 \Gamma_{jyr} \log M}$ 
12       $\delta_2^{(t_2)} = \|\mathbf{A}_j^{(t_2)} - \mathbf{A}_j^{(t_2-1)}\| + \|\mathbf{D}_j^{(t_2)} - \mathbf{D}_j^{(t_2-1)}\| + \|\mathbf{\Gamma}_j^{(t_2)} - \mathbf{\Gamma}_j^{(t_2-1)}\|$ 
13       $t_2 = t_2 + 1$ 
14    end
15     $a_{jr}^{(t_1)*} = a_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})}$ 
16     $\Gamma_{jr}^{(t_1)*} = \Gamma_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})}$ 
17  end
18   $\delta_1^{(t_1)} = \|\mathbf{A}^{(t_1)} - \mathbf{A}^{(t_1-1)}\| + \|\mathbf{D}^{(t_1)} - \mathbf{D}^{(t_1-1)}\| + \|\mathbf{\Gamma}^{(t_1)} - \mathbf{\Gamma}^{(t_1-1)}\|$ 
19   $t_1 = t_1 + 1$ 
20 end
```

Note that $\text{diag}(\hat{\Sigma}_{1r})$ is the r th element on the diagonal of the estimated covariance matrix of the reference group $\hat{\Sigma}_1$.

$\varepsilon_1 = 10^{-3}$ and $\varepsilon_2 = 10^{-7}$.

6.3 Simulation

Sample Size. The total sample size is $N = 3000$, and the group sample sizes are $N_1 = N_2 = N_3 = 1000$.

Test Length. $m = 20$. Simple structure. 10 items per dimension.

Proportion of DIF. 4 items with DIF. 2 DIF items per dimension.

Magnitude of DIF. The first focal group with smaller discrimination parameter (-0.5) on the 4 DIF items.

The second focal group with much smaller difficulty parameter (-1) on the 4 DIF items.

Generated parameters.

$a_{j1} \sim U(1.5, 2.5), j = 1, \dots, 10$

$a_{j2} \sim U(1.5, 2.5), j = 11, \dots, 20$

$$d_1 \sim N(0, 1)$$

$$\mathbf{A} = \begin{pmatrix} 2.17 & 0 \\ 0 & 2.46 \\ 2.41 & 0 \\ 2.45 & 0 \\ 2.34 & 0 \\ 1.84 & 0 \\ 1.85 & 0 \\ 1.92 & 0 \\ 1.94 & 0 \\ 1.90 & 0 \\ 1.92 & 0 \\ 0 & 2.43 \\ 0 & 1.82 \\ 0 & 2.22 \\ 0 & 1.93 \\ 0 & 1.88 \\ 0 & 1.84 \\ 0 & 2.12 \\ 0 & 2.42 \\ 0 & 2.15 \end{pmatrix},$$

$$\mathbf{D} = \begin{pmatrix} 0.03 \\ -1.28 \\ 0.58 \\ -2.06 \\ 0.12 \\ 3.25 \\ -0.41 \\ -0.51 \\ 0.89 \\ 1.33 \\ 0.85 \\ 0.82 \\ -0.37 \\ -0.99 \\ -0.27 \\ 0.19 \\ 1.73 \\ 0.05 \\ -1.86 \\ -0.63 \end{pmatrix},$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \text{ for } j = 1, 2, 3, 6, 7, 8, 9, 10, 11, 14, 15, 16, 17, 18, 19, 20$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} -0.5 & 0 \\ -1 & 0 \end{pmatrix}, \text{ for } j = 4, 5$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & -0.5 \\ 0 & -1 \end{pmatrix}, \text{ for } j = 12, 13$$

No impact. $\theta_i \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.85 \\ 0.85 & 1 \end{pmatrix}\right)$

Tuning parameters. Tuning parameters η are chosen from a sequence starting from 21 to 51 with increment 3.

6.3.1 Results of 20 Replications

Table 7. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.9583	0.2916	0.9583
Type I	0.02778	0.0238	0.0119

Omnibus DIF is defined as if at least one focal group showd DIF on an item, then that item is flagged as DIF.

Table 8. Type I error and Power of mirt LRT

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.9027	0	0.9305
Type I	0.02777	0.00396	0.01587

Table 9. Type I error and Power of mirt Wald

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0	0	0.9305
Type I	0.00396	0.00396	0.0119

Both regularization and mirt LRT can detect DIF maginitude 1 with power 100%. Our regularization method has slightly lower Type I error.

Table 10. Item parameter estimates by regularization

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0042	-0.008	-0.0100
RMSE	0.1810	0.1554	0.1019

Table 11. Item parameter estimates by mirt LRT (three groups togrther)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0241	0.0065	-0.0056
RMSE	0.1627	0.1457	0.09548

Item parameter estimates by mirt LRT (reference and focal 1)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0235	0.0136	-0.0114

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
RMSE	0.1821	0.1643	0.1157

Item parameter estimates by mirt LRT (reference and focal 2)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0239	0.0121	-0.0098
RMSE	0.1817	0.1647	0.1089

Table 12. Item parameter estimates by mirt Wald (three groups together)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0247	0.0109	-0.0073
RMSE	0.1814	0.1692	0.0958

Item parameter estimates by mirt LRT (reference and focal 1)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0235	0.0136	-0.0110
RMSE	0.1830	0.1651	0.1157

Item parameter estimates by mirt LRT (reference and focal 2)

Item Parameters	\mathbf{a}_1	\mathbf{a}_2	\mathbf{d}
Bias	0.0239	0.0121	-0.0098
RMSE	0.1817	0.1647	0.1089

LRT method has slightly better non-DIF item parameter estimates when including three groups.

Table 13. Absolute bias for DIF magnitude recoveries that were true DIF

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Regularization	0.343	0.396	0.290
mirt LRT (three groups together)	0.200	0.195	0.204
mirt LRT (reference and focal 1)	-	0.2818	-
mirt LRT (reference and focal 2)	-	-	0.2157
mirt Wald (three groups together)	0.3453	0.2647	0.4259
mirt Wald (reference and focal 1)	-	0.2617	-
mirt Wald (reference and focal 2)	-	-	0.2157

Absolute bias for 4 items (4,5,12,13) with DIF (include false negative).

One advantage compared to mirt multipleGroup is that our method (re-fit part) allows some of focal groups to have DIF while others do not.

One drawback of our method is when estimating three groups together, low power in the first focal group would affect the parameter estimate of the second focal group.

7 Non-uniform DIF Detection via Group LASSO

When the items have non-uniform DIF on both slope and intercept, the DIF parameter we are estimating are $\mathbf{\Gamma} = (\mathbf{0}, \dots, \mathbf{0}, \mathbf{\Gamma}_{q+1}, \dots, \mathbf{\Gamma}_m)$ and $\boldsymbol{\beta} = (\mathbf{0}, \dots, \mathbf{0}, \boldsymbol{\beta}_{q+1}, \dots, \boldsymbol{\beta}_m)$.

7.1 E step

The equations for expectations are same as before.

In our DIF detection problem, we minimize the following objective function

$$-\log M + \eta \sum_j^m \|(\mathbf{\Gamma}_j, \boldsymbol{\beta}_j)\|_2 \quad (52)$$

For each item, we minimize

$$-\log M_j + \eta \|(\mathbf{\Gamma}_j, \boldsymbol{\beta}_j)\|_2 \quad (53)$$

where η is the group lasso tuning parameter.

We denote by $\boldsymbol{\tau} \in \mathbb{R}^{(y-1)*r+(y-1)*(m-1)}$ the whole DIF parameter vector, i.e. $\boldsymbol{\tau} = (\mathbf{\Gamma}, \boldsymbol{\beta})^T$.

Then, our objective function is

$$S_\eta(\boldsymbol{\tau}) = -\log M + \eta \sum_j^m \|\boldsymbol{\tau}\|_2. \quad (54)$$

For each item j ,

$$S_\eta(\boldsymbol{\tau}_j) = -\log M_j + \eta \|\boldsymbol{\tau}_j\|_2. \quad (55)$$

7.2 M step

The equations are same as in section 3.2. The Block co-ordinate gradient descent for solving the group lasso problem as follow.

7.2.1 Block co-ordinate gradient descent

Using a second-order Taylor series expansion at $\hat{\boldsymbol{\tau}}^{(t-1)}$ we define

$$M_\eta^{(t-1)}(\boldsymbol{\epsilon}^{(t)}) = -\{\log M + \boldsymbol{\epsilon}^{(t)T} \nabla \log M + \frac{1}{2} \boldsymbol{\epsilon}^{(t)T} H^{(t-1)} \boldsymbol{\epsilon}^{(t)}\} + \eta \sum_j^m \|\hat{\boldsymbol{\tau}}_j^{(t-1)} + \boldsymbol{\epsilon}_j^{(t)}\|_2,$$

where $\boldsymbol{\tau}^{(t)} = \boldsymbol{\tau}^{(t-1)} + \boldsymbol{\epsilon}^{(t)}$, and

$$\nabla \log M = \left(\frac{\partial \log M}{\partial \mathbf{\Gamma}}, \frac{\partial \log M}{\partial \boldsymbol{\beta}} \right)$$

and

$$H^{(t-1)} = \begin{pmatrix} \frac{\partial^2 \log M}{\partial \Gamma^2} & \frac{\partial^2 \log M}{\partial \Gamma \partial \beta} \\ \frac{\partial^2 \log M}{\partial \Gamma \partial \beta} & \frac{\partial^2 \log M}{\partial \beta^2} \end{pmatrix}.$$

We have $M_\eta^{(t-1)}(\epsilon) \approx S_\eta(\hat{\tau}^{(t-1)} + \epsilon^{(t)})$.

We run a block co-ordinate gradient descent algorithm for each group (item) with all other groups fixed.

For item j , denote u to be the subgradient of $\|\tau_j^{(t-1)} + \epsilon_j^{(t)}\|_2$. We have

$$u = \begin{cases} \frac{\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}}{\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2}, & \text{if } \hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)} \neq \mathbf{0} \\ \in \{u : \|u\|_2 \leq 1\}, & \text{if } \hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)} = \mathbf{0} \end{cases}.$$

The subgradient equation $\partial_{\epsilon_j} M_\eta^{(t-1)}(\epsilon^{(t)}) = -\nabla \log M_j - \epsilon_j^{(t)T} H_j^{(t-1)} + \eta u = 0$ is satisfied with $\tau_j^{(t-1)} + \epsilon_j = 0$ if

$$\|u\|_2 = \left\| \frac{\nabla \log M_j + \epsilon_j^{(t)T} H_j^{(t-1)}}{\eta} \right\|_2 \leq 1$$

$$\|\nabla \log M_j + \epsilon_j^{(t)T} H_j^{(t-1)}\|_2 \leq \eta$$

$$\|\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}\|_2 \leq \eta,$$

the minimizer of $M_\eta^{(t-1)}(\epsilon)$ is

$$\hat{\epsilon}_j^{(t)} = -\hat{\tau}_j^{(t-1)}.$$

Otherwise,

Then the subgradient equation is

$$\partial_{\epsilon_j} M_\eta^{(t-1)}(\epsilon^{(t)}) = -\nabla \log M_j - \epsilon_j^{(t)T} H_j^{(t-1)} + \eta \frac{\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}}{\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2} = 0 \quad (56)$$

$$\partial_{\epsilon_j} M_\eta^{(t-1)}(\epsilon^{(t)}) = -\nabla \log M_j - (\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}) H_j^{(t-1)} + \hat{\tau}_j^{(t-1)} H_j^{(t-1)} + \eta \frac{\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}}{\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2} = 0$$

$$\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)} = -(\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}) H_j^{(t-1)} + \eta \frac{\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}}{\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2} \quad (57)$$

$$\tau = \hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)} = \frac{(\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}) \|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2}{\eta - H_j^{(t-1)} \|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2} \quad (58)$$

Taking the norm of both sides of (57) we see that

$$\|\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}\|_2 = \left(\frac{\eta}{\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2} - H_j^{(t-1)} \right) \|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2$$

$$\|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2 = \frac{\eta - \|\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}\|_2}{H_j^{(t-1)}} \quad (59)$$

Plugging (59) into (58), we have

$$\epsilon_j^{(t)} = -(H_j^{(t-1)})^{-1} \left\{ \nabla \log M_j - \eta \frac{\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}}{\|\nabla \log M_j - \hat{\tau}_j^{(t-1)} H_j^{(t-1)}\|_2} \right\}. \quad (60)$$

$$\nabla \log M_j = \left(\frac{\partial \log M}{\partial \gamma_{j11}}, \dots, \frac{\partial \log M}{\partial \gamma_{jyr}}, \dots, \frac{\partial \log M}{\partial \gamma_{j3q}}, \frac{\partial \log M}{\partial \beta_{j11}}, \dots, \frac{\partial \log M}{\partial \beta_{jyk}}, \dots, \frac{\partial \log M}{\partial \beta_{j3(p-1)}} \right), r = 1, \dots, q; k = 1, \dots, p-1; y = 2, 3.$$

If $\epsilon_j^{(t)} \neq 0$, performing a Backtracking-Armijo line search: let $\alpha^{(t)}$ be the largest value in $\{\alpha_0 \delta^l\}_{l \geq 0}$ s.t.

$$S_\eta(\hat{\tau}_j^{(t-1)} + \alpha^{(t)} \epsilon_j^{(t)}) - S_\eta(\hat{\tau}_j^{(t-1)}) \leq \alpha^{(t)} \sigma \Delta^{(t)}, \quad (61)$$

where $\alpha_0 = 1$, $\delta = 0.5$ and $\sigma = 0.1$, and Δ is the improvement in the objective function $S_\eta(\cdot)$ when using a linear approximation for the log-likelihood, i.e.

$$\Delta_j^{(t)} = -\epsilon_j^{(t)T} \nabla \log M_j + \eta \|\hat{\tau}_j^{(t-1)} + \epsilon_j^{(t)}\|_2 - \eta \|\hat{\tau}_j^{(t-1)}\|_2. \quad (62)$$

$$\hat{\tau}_j^{(t)} = (\hat{\Gamma}_j^{(t)}, \hat{\beta}_j^{(t)}) = \hat{\tau}_j^{(t-1)} + \alpha^{(t)} \epsilon_j^{(t)}$$

Then we update $P_{jky|q_a}^*$ and $Q_{jky|q_a}^*$ by plugging in $\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\Gamma}$ and $\hat{\beta}$ from last coordinate descent cycle and repeat above steps until a convergence criterion is met.

After we get optimizers for item j , we do transformations on all estimates as following

$$a_{jr}^{(t)*} = a_{jr}^{(t)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})},$$

$$\gamma_{jr}^{(t)*} = \gamma_{jr}^{(t)} * \sqrt{\text{diag}(\hat{\Sigma}_{1r})},$$

where μ_{1r} is the r th element of the estimated mean vector of the reference group $\hat{\mu}_1$, and $\text{diag}(\hat{\Sigma}_{1r})$ is the r th element on the diagonal of the estimated covariance matrix of the reference group $\hat{\Sigma}_1$.

We run an EM and block coordinate gradient descent algorithm given by following.

Algorithm 3: Non-uniform DIF Detection via group LASSO

Input : $\mathbf{A}_0, \mathbf{D}_0, \mathbf{\Gamma}_0, \boldsymbol{\beta}_0, \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0, \mathbf{U}, \eta, \varepsilon_1, \varepsilon_2$

Output: $\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\mathbf{\Gamma}}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}$

```

1 set  $t_1 = 1$ ,  $\delta^{(t_1-1)} = \text{any value greater than } \varepsilon_1$ 
2 while  $\delta_1^{(t_1-1)} > \varepsilon_1$  do
3   Calculate  $n_{gy}$  and  $r_{gjk}$ 
4   Update  $\boldsymbol{\mu}^{(t_1)}$  and  $\boldsymbol{\Sigma}^{(t_1)}$ 
5   for  $j=1, \dots, m$  do
6     set  $t_2 = 1$ ,  $\delta_2^{(t_2-1)} = \text{any value greater than } \varepsilon_2$ 
7     while  $\delta_2^{(t_2-1)} > \varepsilon_2$  do
8       Calculate  $P_{jyk|q_g}^*, Q_{jyk|q_g}^*$ 
9        $a_{jr}^{(t_2)} = a_{jr}^{(t_2-1)} - \frac{\partial_{a_{jr}} \log M}{\partial_{a_{jr}}^2 \log M}$ 
10       $d_{jk}^{(t_2)} = d_{jk}^{(t_2-1)} - \frac{\partial_{d_{jk}} \log M}{\partial_{d_{jk}}^2 \log M}$ 
11      if  $\|\nabla \log M_j - \hat{\boldsymbol{\tau}}_j^{(t_2-1)} H_j^{(t_2-1)}\|_2 \leq \eta$  then
12         $\boldsymbol{\tau}_j = \mathbf{0}$ 
13      else
14         $\boldsymbol{\epsilon}_j^{(t_2)} = -(H_j^{(t_2-1)})^{-1} \{ \nabla \log M_j - \eta \frac{\nabla \log M_j - \hat{\boldsymbol{\tau}}_j^{(t_2-1)} H_j^{(t_2-1)}}{\|\nabla \log M_j - \hat{\boldsymbol{\tau}}_j^{(t_2-1)} H_j^{(t_2-1)}\|_2} \}$ 
15         $\Delta_j^{(t_2)} = -\boldsymbol{\epsilon}_j^{(t_2)T} \nabla \log M_j + \eta \|\hat{\boldsymbol{\tau}}_j^{(t_2-1)} + \boldsymbol{\epsilon}_j^{(t_2)}\|_2 - \eta \|\hat{\boldsymbol{\tau}}_j^{(t_2-1)}\|_2$ 
16         $\alpha^{(t_2)}$  is the max value in  $\{\alpha^{(0)} \delta^l\}_{l \geq 0}$  such that
17         $S_\eta(\hat{\boldsymbol{\tau}}_j^{(t_2-1)} + \alpha^{(t_2)} \boldsymbol{\epsilon}_j^{(t_2)}) - S_\eta(\hat{\boldsymbol{\tau}}_j^{(t_2-1)}) \leq \alpha^{(t_2)} \sigma \Delta^{(t_2)}$ .
18         $\boldsymbol{\tau}_j^{(t_2)} = \boldsymbol{\tau}_j^{(t_2-1)} + \alpha^{(t_2)} \boldsymbol{\epsilon}_j^{(t_2)}$ 
19      end
20       $\delta_2^{(t_2)} = \|\mathbf{A}_j^{(t_2)} - \mathbf{A}_j^{(t_2-1)}\| + \|\mathbf{D}_j^{(t_2)} - \mathbf{D}_j^{(t_2-1)}\| + \|\alpha^{(t_2-1)} \boldsymbol{\epsilon}_j^{(t_2-1)}\|$ 
21       $t_2 = t_2 + 1$ 
22    end
23     $a_{jr}^{(t_1)*} = a_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})}$ 
24     $\Gamma_{jr}^{(t_1)*} = \Gamma_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})}$ 
25  end
26   $\delta_1^{(t_1)} = \|\mathbf{A}^{(t_1)} - \mathbf{A}^{(t_1-1)}\| + \|\mathbf{D}^{(t_1)} - \mathbf{D}^{(t_1-1)}\| + \|\mathbf{\Gamma}^{(t_1)} - \mathbf{\Gamma}^{(t_1-1)}\| + \|\boldsymbol{\beta}^{(t_1)} - \boldsymbol{\beta}^{(t_1-1)}\|$ 
27   $t_1 = t_1 + 1$ 
28 end
```

Note that $\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})$ is the r th element on the diagonal of the estimated covariance matrix of the reference group $\hat{\boldsymbol{\Sigma}}_1$.

$\varepsilon_1 = 10^{-3}$ and $\varepsilon_2 = 10^{-7}$.

7.3 Simulation

Sample Size. The total sample size is $N = 3000$, and the group sample sizes are $N_1 = N_2 = N_3 = 1000$.

Test Length. $m = 20$. Simple structure. 10 items per dimension.

Proportion of DIF. 4 items with DIF. 2 DIF items per dimension.

Magnitude of DIF. The first focal group with smaller discrimination parameter (-0.5) and larger difficulty parameters (+0.5) on the 4 DIF items.

The second focal group with much smaller difficulty parameter (-1) and much larger difficulty parameters (+1) on the 4 DIF items.

Generated parameters.

$$a_{j1} \sim U(1.5, 2.5), j = 1, \dots, 10$$

$$a_{j2} \sim U(1.5, 2.5), j = 11, \dots, 20$$

$$d_1 \sim N(0, 1)$$

$$\mathbf{A} = \begin{pmatrix} 2.17 & 0 \\ 0 & 2.46 \\ 2.41 & 0 \\ 2.45 & 0 \\ 2.34 & 0 \\ 1.84 & 0 \\ 1.85 & 0 \\ 1.92 & 0 \\ 1.94 & 0 \\ 1.90 & 0 \\ 1.92 & 0 \\ 0 & 2.43 \\ 0 & 1.82 \\ 0 & 2.22 \\ 0 & 1.93 \\ 0 & 1.88 \\ 0 & 1.84 \\ 0 & 2.12 \\ 0 & 2.42 \\ 0 & 2.15 \end{pmatrix},$$

$$\mathbf{D} = \begin{pmatrix} 0.03 \\ -1.28 \\ 0.58 \\ -2.06 \\ 0.12 \\ 3.25 \\ -0.41 \\ -0.51 \\ 0.89 \\ 1.33 \\ 0.85 \\ 0.82 \\ -0.37 \\ -0.99 \\ -0.27 \\ 0.19 \\ 1.73 \\ 0.05 \\ -1.86 \\ -0.63 \end{pmatrix},$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \text{ for } j = 1, 2, 3, 6, 7, 8, 9, 10, 11, 14, 15, 16, 17, 18, 19, 20$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} -0.3 & 0 \\ -0.6 & 0 \end{pmatrix}, \text{ for } j = 4, 5$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & -0.3 \\ 0 & -0.6 \end{pmatrix}, \text{ for } j = 12, 13$$

$$\boldsymbol{\beta} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.5 & 1 \\ 0.5 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0.5 & 1 \\ 0.5 & 1 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix},$$

No impact. $\theta_i \sim N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0.85 \\ 0.85 & 1 \end{pmatrix}\right)$

Tuning parameters. Tuning parameters η are chosen from a sequence starting from 21 to 51 with increment 3.

7.3.1 Results of 20 Replications

Table 14. *Type I error and Power of regularization method*

Group	Omnibus DIF
Power	0.975
Type I	0.065

Omnibus DIF is defined as if at least one focal group showed DIF on an item, then that item is flagged as DIF.