

# DIF Detection via Regularization

Ruoyi Zhu

9/9/2020

## Contents

<b>1</b>	<b>Conventional DIF Detection Methods</b>	<b>1</b>
1.1	Likelihood Ratio Test . . . . .	1
<b>2</b>	<b>Multidimensional 2PL Model with DIF</b>	<b>2</b>
<b>3</b>	<b>Graded Response Model with DIF</b>	<b>2</b>
<b>4</b>	<b>Model Identifiability Constraint</b>	<b>5</b>
<b>5</b>	<b>Uniform DIF Detection via LASSO</b>	<b>6</b>
5.1	E step . . . . .	6
5.2	M step . . . . .	8
5.3	Tuning Parameter Selection . . . . .	13
5.4	Simulation . . . . .	13
<b>6</b>	<b>Non-uniform DIF Detection via LASSO</b>	<b>18</b>
6.1	E step . . . . .	19
6.2	M step . . . . .	19
6.3	Tuning Parameter Selection . . . . .	21
6.4	Simulation . . . . .	21
<b>7</b>	<b>Non-uniform DIF Detection via Group LASSO</b>	<b>25</b>
7.1	E step . . . . .	25
7.2	M step . . . . .	26
7.3	Tuning Parameter Selection . . . . .	29
7.4	Simulation . . . . .	30
<b>8</b>	<b>Plots of Results</b>	<b>35</b>
8.1	DIF only on intercept . . . . .	35
8.2	DIF only on slope . . . . .	36
8.3	DIF on slope and intercept . . . . .	37
<b>9</b>	<b>High DIF Proportion Issue</b>	<b>37</b>

## 1 Conventional 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  $\boldsymbol{\theta}_i$  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.

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

is the non-uniform DIF parameter, where  $\boldsymbol{\gamma}_{j1}$  is the non-uniform DIF parameter for the first focal group and  $\boldsymbol{\gamma}_{j2}$  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  $\beta_{j1}$  is the uniform DIF parameter for the first focal group and  $\beta_{j2}$  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

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

boundary parameter

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

non-uniform DIF parameter

$$\boldsymbol{\Gamma} = (\boldsymbol{\gamma}_1, \boldsymbol{\gamma}_2, \dots, \boldsymbol{\gamma}_j, \dots, \boldsymbol{\gamma}_m) \quad (j = 1, \dots, m)$$

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

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

and uniform DIF parameter

$$\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_j, \dots, \boldsymbol{\beta}_m) \quad (j = 1, \dots, m)$$

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

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

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

The  $N * m$  response matrix is

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

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 rfirst 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

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

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  $f(\mathbf{q}_g \mid \boldsymbol{\mu}_y^{(t-1)}, \boldsymbol{\Sigma}_y^{(t-1)})$  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

Some constraints on the item parameters are required for model identification. Here, for each dimension, we set one anchor item which we know its DIF parameters ( $\Gamma$  and  $\beta$ ) 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

$$\beta = (\mathbf{0}, \mathbf{0}, \beta_3, \dots, \beta_m).$$

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  $\Gamma = \mathbf{0}$ . The DIF parameter we are estimating is only  $\beta = (\mathbf{0}, \dots, \mathbf{0}, \beta_{q+1}, \dots, \beta_m)$ .

### 5.1 E step

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

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

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

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

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

$$m(\mathbf{A}, \mathbf{D}, \beta, \mu_y, \Sigma_y \mid \mathbf{y}, \mathbf{u}_i) = \int L(\mathbf{A}, \mathbf{D}, \beta \mid \mathbf{y}, \mathbf{u}_i, \theta_i) f(\mu_y, \Sigma_y \mid \theta_i) \partial \theta_i \quad (6)$$

Then

$$h(\theta_i \mid \mathbf{u}_i, \mathbf{y}_i, \mathbf{A}^{(t-1)}, \mathbf{D}^{(t-1)}, \beta^{(t-1)}, \mu_y^{(t-1)}, \Sigma_y^{(t-1)}) = \frac{L(\mathbf{A}, \mathbf{D}, \beta \mid \mathbf{y}, \mathbf{u}_i, \theta_i) f(\mu_y, \Sigma_y \mid \theta_i)}{m(\mathbf{A}, \mathbf{D}, \beta, \mu_y, \Sigma_y \mid \mathbf{y}, \mathbf{u}_i)} \quad (7)$$

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

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

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

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  $\mathbf{q}_g$ .

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

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  $\eta$  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  $\mu$  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  $\beta_{jyk}$ , 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

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

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

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

and

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

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

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

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

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

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  $\tau$ , 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  $\beta$ , 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 Tuning Parameter Selection

The Bayesian information criterion (BIC) is applied for tuning parameter selection. The BIC can be calculated by

$$BIC_{\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\beta}} = -2 \max_{\mathbf{A}, \mathbf{D}, \beta} \log M + \|\hat{\beta}\|_0 \log N, \quad (11)$$

where  $N$  is the sample size. The first term control the bias of the estimator, and the second term penalize the complexity of the model. We want to choose a tuning parameter  $\eta$  which gives us the smallest BIC value.

### 5.4 Simulation

#### 5.4.1 Simulation 1

*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.*

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

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

$d_1 \sim N(0, 1)$

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

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

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

wABC

item	4	5	12	13
Focal 1	0.886	1.052	0.998	1.275
Focal 2	1.857	2.068	1.924	2.537

The first two items are used as anchor items in estimation.

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  $\eta$  are chosen from a sequence starting from 21 to 51 with increment 3. The range [21,51] was chosen because in all simulation studies I've run so far, the true model will be selected by the tuning parameter in this range for 100%.

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 in regularization is defined as if at least one focal group showed DIF on an item, then that item is flagged as DIF. But in mirt package, omnibus DIF is testing  $H_0$ : the item has no DIF v.s.  $H_\alpha$ : all focal groups have 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 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 for true DIF items*

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. Trait Estimation

Item Parameters	$\theta_1$	$\theta_2$
Bias	0.0365	0.0374
RMSE	0.3823682	0.3806336
SE	0.3644813	0.3590251

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.

#### 5.4.2 Simulation 2

Increase DIF proportion to 60% and keep everything else the same in simulation 1.

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

*Proportion of DIF.* 12 items with DIF. 6 DIF items per dimension.

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

wABC

item	4	5	6	7	8	9	12	13	14	15	16	17
Focal 1	0.886	1.052	0.636	1.265	1.255	1.116	0.998	1.275	1.091	1.234	1.233	0.918
Focal 2	1.857	2.067	1.385	2.521	2.513	2.117	1.923	2.537	2.218	2.448	2.399	1.684

Table 7. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.872	0.136	0.872
Type I	0.135	0.085	0.053



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 (0.05 significance level)

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.963	0.395	0.987
Type I	0.0175	0.0125	0.0175

### 5.4.3 Simulation 3

Increase total sample size to 3000 and keep everything else the same in simulation 1.

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

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

Table 9. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	1	0.775	1
Type I	0.051	0.027	0.028

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

Table 10. 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	0.995	0.63	0.995
Type I	0.008	0.01	0.012

Table 11. Item parameter estimates by regularization

Item Parameters	$\mathbf{a}_1$	$\mathbf{a}_2$	$\mathbf{d}$
Bias	-0.01364	-0.00353	-0.0316
RMSE	0.11486	0.11917	0.111221

Table 12. Item parameter estimates by mirt LRT (0.05 significance level)

Item Parameters	$\mathbf{a}_1$	$\mathbf{a}_2$	$\mathbf{d}$
Bias	0.00612	0.01542	0.009
RMSE	0.11613	0.12292	0.0953

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 (include false negative)	0.1111	0.1884	0.1451
mirt LRT (include false negative)	0.1145	0.2557	0.1155

#### 5.4.4 Simulation 4

Increase total sample size to 3000 and keep everything else the same in simulation 2.

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

*Proportion of DIF.* 12 items with DIF. 6 DIF items per dimension.

Table 14. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.998	0.58	0.998
Type I	0.05	0.04	0.01

Table 15. 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	0.996	0.772	1
Type I	0.02	0.02	0.0225

Table 16. Item parameter estimates by regularization

Item Parameters	$\mathbf{a}_1$	$\mathbf{a}_2$	$\mathbf{d}$
Bias	-0.0316	0.0244	0.0192
RMSE	0.118	0.127	0.131

Table 17. Item parameter estimates by mirt LRT (0.05 significance level)

Item Parameters	$\mathbf{a}_1$	$\mathbf{a}_2$	$\mathbf{d}$
Bias	0.0075	0.0197	0.0093
RMSE	0.1188	0.1269	0.1062

Table 18. 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)	0.216	0.2663	0.167
mirt LRT (include false negative)	0.1843	0.3556	0.1155

## 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  $\beta$  by  $\Gamma$ , and minimize the following objective function

$$-\log M + \eta \sum_j^m ||\Gamma_j||_1 \quad (43)$$

For each item, we minimize

$$-\log M_j + \eta ||\Gamma_j||_1 \quad (44)$$

where  $\eta$  is the lasso tuning parameter.

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

## 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  $\gamma_{jry}$ , 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).

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

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

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

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

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

### 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

$$\begin{aligned} a_{jr}^{(t)} &= a_{jr}^{(t-1)} - \frac{\partial_{a_{jr}} \log M}{\partial_{a_{jr}}^2 \log M}, \\ d_{jk}^{(t)} &= d_{jk}^{(t-1)} - \frac{\partial_{d_{jk}} \log M}{\partial_{d_{jk}}^2 \log M} \end{aligned}$$

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  $\hat{\mathbf{\beta}}$  from last coordinate descent cycle and repeat above steps until a convergence criterion is met.

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

$$\begin{aligned} a_{jr}^{(t)*} &= a_{jr}^{(t)} * \sqrt{\operatorname{diag}(\hat{\mathbf{\Sigma}}_{1r})}, \\ \gamma_{jr}^{(t)*} &= \gamma_{jr}^{(t)} * \sqrt{\operatorname{diag}(\hat{\mathbf{\Sigma}}_{1r})}, \end{aligned}$$

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

---

**Algorithm 2:** Non-uniform DIF Detection via LASSO

---

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

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

---

```

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  $r_{ggyk}$ 
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       $\Gamma_{jyr}^{(t_2)} = -\frac{\text{soft}(\partial_{\Gamma_{jyr}} \log M - \Gamma_{jyr}^{(t_2-1)} * \partial_{\Gamma_{jyr}}^2 \log M, \eta)}{\partial_{\Gamma_{jyr}}^2 \log M}$ 
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     $a_{jr}^{(t_1)*} = a_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})}$ 
16     $\Gamma_{jr}^{(t_1)*} = \Gamma_{jr}^{(t_1)} * \sqrt{\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})}$ 
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  $\text{diag}(\hat{\boldsymbol{\Sigma}}_{1r})$  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  $\varepsilon_2 = 10^{-7}$ .

### 6.3 Tuning Parameter Selection

The BIC can be calculated by

$$BIC_{\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\mathbf{\Gamma}}} = -2 \max_{\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\mathbf{\Gamma}}} \log M + \|\hat{\mathbf{\Gamma}}\|_0 \log N, \quad (11)$$

where  $N$  is the sample size.

### 6.4 Simulation

#### 6.4.1 Simulation 5

*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 smaller discrimination parameter (-0.5) on the 4 DIF items.  
The second focal group with much smaller discrimination parameter (-1) on the 4 DIF items.

*Generated parameters.*

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

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

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

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

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

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

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

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

wABC

item	4	5	12	13
Focal 1	0.620	0.530	0.560	0.774
Focal 2	1.331	1.242	1.292	1.813

The first two items are used as anchor items in estimation.

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  $\eta$  are chosen from a sequence starting from 21 to 51 with increment 3.

Results of 50 Replications

Table 19. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.625	0.14	0.625
Type I	0.0314	0.0171	0.0185

Table 20. Type I error and Power of mirt LRT

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.515	0.375	0.646
Type I	0.0185	0.0238	0.0149

### 6.4.2 Simulation 6

Increase DIF proportion to 60%, and keep everything else the same as simulation 5.

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

*Proportion of DIF.* 12 items with DIF. 6 DIF items per dimension.

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \text{ for } j = 1, 2, 3, 10, 11, 18, 19, 20$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} -0.5 & 0 \\ -1 & 0 \end{pmatrix}, \text{ for } j = 4, 5, 6, 7, 8, 9$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & -0.5 \\ 0 & -1 \end{pmatrix}, \text{ for } j = 12, 13, 14, 15, 16, 17$$

wABC

item	4	5	6	7	8	9	12	13	14	15	16	17
Focal 1	0.620	0.530	0.556	0.764	0.760	0.737	0.560	0.774	0.635	0.724	0.744	0.750
Focal 2	1.331	1.242	1.076	1.788	1.771	1.686	1.292	1.813	1.452	1.702	1.751	1.589

Table 21. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.45	0.06	0.443
Type I	0.0067	0.0033	0.0033

Adaptive LASSO (selected by GIC)

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.383	0.0086	0.3836
Type I	0	0	0

Table 22. Type I error and Power of mirt LRT

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.656	0.185	0.746
Type I	0.0067	0.01	0.0233

### 6.4.3 Simulation 7

Increase total sample size to 3000, and keep everything else the same as simulation 5.

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

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

Table 23. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.96	0.345	0.96
Type I	0.0371	0.0271	0.0214

Table 24. Type I error and Power of mirt LRT

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.815	0.555	0.885
Type I	0.005	0.031	0.0119

### 6.4.4 Simulation 8

Increase DIF proportion to 60%, and keep everything else the same as simulation 7.



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

*Proportion of DIF.* 12 items with DIF. 6 DIF items per dimension.

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \text{ for } j = 1, 2, 3, 10, 11, 18, 19, 20$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} -0.5 & 0 \\ -1 & 0 \end{pmatrix}, \text{ for } j = 4, 5, 6, 7, 8, 9$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & -0.5 \\ 0 & -1 \end{pmatrix}, \text{ for } j = 12, 13, 14, 15, 16, 17$$

Table 25. Type I error and Power of regularization method

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.9933	0.481	0.9933
Type I	0.0333	0.02	0.02

Adaptive LASSO (selected by GIC)

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.866	0.112	0.866
Type I	0.0083	0.0083	0

Table 26. Type I error and Power of mirt LRT

Group	Omnibus DIF	Group with DIF=0.5	Group with DIF=1
Power	0.9367	0.3	0.9683
Type I	0.0233	0.0133	0.0333

## 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  $\eta$  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

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

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

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

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

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  $\|\boldsymbol{\tau}_j^{(t-1)} + \boldsymbol{\epsilon}_j^{(t)}\|_2$ . We have

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

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

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

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

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

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

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

Otherwise,

Then the subgradient equation is

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

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

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

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

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

$$\begin{aligned} \|\nabla \log M_j - \hat{\boldsymbol{\tau}}_j^{(t-1)} H_j^{(t-1)}\|_2 &= \left( \frac{\eta}{\|\hat{\boldsymbol{\tau}}_j^{(t-1)} + \boldsymbol{\epsilon}_j^{(t)}\|_2} - H_j^{(t-1)} \right) \|\hat{\boldsymbol{\tau}}_j^{(t-1)} + \boldsymbol{\epsilon}_j^{(t)}\|_2 \\ \|\hat{\boldsymbol{\tau}}_j^{(t-1)} + \boldsymbol{\epsilon}_j^{(t)}\|_2 &= \frac{\eta - \|\nabla \log M_j - \hat{\boldsymbol{\tau}}_j^{(t-1)} H_j^{(t-1)}\|_2}{H_j^{(t-1)}} \end{aligned} \quad (59)$$

Plugging (59) into (58), we have

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

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

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

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

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

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

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

Then we update  $P_{jky|q_g}^*$  and  $Q_{jky|q_g}^*$  by plugging in  $\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\boldsymbol{\Gamma}}$  and  $\hat{\boldsymbol{\beta}}$  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

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

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

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

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

---

```

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  $r_{ggyk}$ 
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 Tuning Parameter Selection

The Bayesian information criterion (BIC) in Tutz et al.(2015) is applied for tuning parameter selection. The degrees of freedom of penalized parameters  $\boldsymbol{\tau} = (\boldsymbol{\Gamma}, \boldsymbol{\beta})^T$  are approximated by

$$\tilde{df}_{\boldsymbol{\tau}}(\eta) = \sum_j^m I(\|\boldsymbol{\tau}_j(\eta)\|_2 > 0) + \sum_j^m \frac{\|\boldsymbol{\tau}_j(\eta)\|_2}{\|\boldsymbol{\tau}_j^{ML}\|_2} (m' - 1)$$

where  $m'$  is the number of DIF parameters for each item. Here we have  $m' = 4$ .

The total degrees of freedom are

$$df(\eta) = m + N + \tilde{df}_{\boldsymbol{\tau}}(\eta) - 1.$$

The BIC can be calculated by

$$BIC(\eta) = -2 \max_{\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\boldsymbol{\tau}}} \log M + df(\eta) \log(m \cdot N).$$

where  $N$  is the sample size. The first term control the bias of the estimator, and the second term penalize the complexity of the model. We want to choose a tuning parameter  $\eta$  which gives us the smallest BIC value.

## 7.4 Simulation

### 7.4.1 Simulation 9

*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 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.*

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

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

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

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

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

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

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

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

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

The first two items are used as anchor items in estimation.

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  $\eta$  are chosen from a sequence starting from 21 to 51 with increment 3.

Results of 50 Replications

Table 27. Type I error and Power of regularization method

Group	Omnibus DIF
Power	0.93
Type I	0.0457

Table 28. Type I error and Power of mirt LRT

Group	Omnibus DIF
Power	0.925
Type I	0.0142

### 7.4.2 Simulation 10

Keep everything the same as simulation 9. Increase the DIF proportion to 60%. *Sample Size.* The total sample size is  $N = 1500$ , and the group sample sizes are  $N_1 = N_2 = N_3 = 500$ .

*Proportion of DIF.* 12 items with DIF. 6 DIF items per dimension.

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \text{ for } j = 1, 2, 3, 10, 11, 18, 19, 20$$



$$\mathbf{\Gamma}_j = \begin{pmatrix} -0.3 & 0 \\ -0.6 & 0 \end{pmatrix}, \text{ for } j = 4, 5, 6, 7, 8, 9$$

$$\mathbf{\Gamma}_j = \begin{pmatrix} 0 & -0.3 \\ 0 & -0.6 \end{pmatrix}, \text{ for } j = 12, 13, 14, 15, 16, 17$$

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

Results of 50 Replications

Table 29. Type I error and Power of regularization method

Group	Omnibus DIF
Power	0.842
Type I	0.0033

Table 30. Type I error and Power of mirt LRT

Group	Omnibus DIF
Power	0.926
Type I	0.023

### 7.4.3 Simulation 11

Keep everything the same as simulation 9. Increase the sample size. *Sample Size*. The total sample size is  $N = 3000$ , and the group sample sizes are  $N_1 = N_2 = N_3 = 1000$ .

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

Results of 50 Replications

Table 31. Type I error and Power of regularization method

Group	Omnibus DIF
Power	1
Type I	0.0557

Table 32. Type I error and Power of mirt LRT

Group	Omnibus DIF
Power	0.985
Type I	0.008

#### 7.4.4 Simulation 12

Keep everything the same as simulation 10. Increase the sample size to 3000. *Sample Size*. The total sample size is  $N = 3000$ , and the group sample sizes are  $N_1 = N_2 = N_3 = 1000$ .

*Proportion of DIF*. 12 items with DIF. 6 DIF items per dimension.

Results of 50 Replications

Table 33. Type I error and Power of regularization method

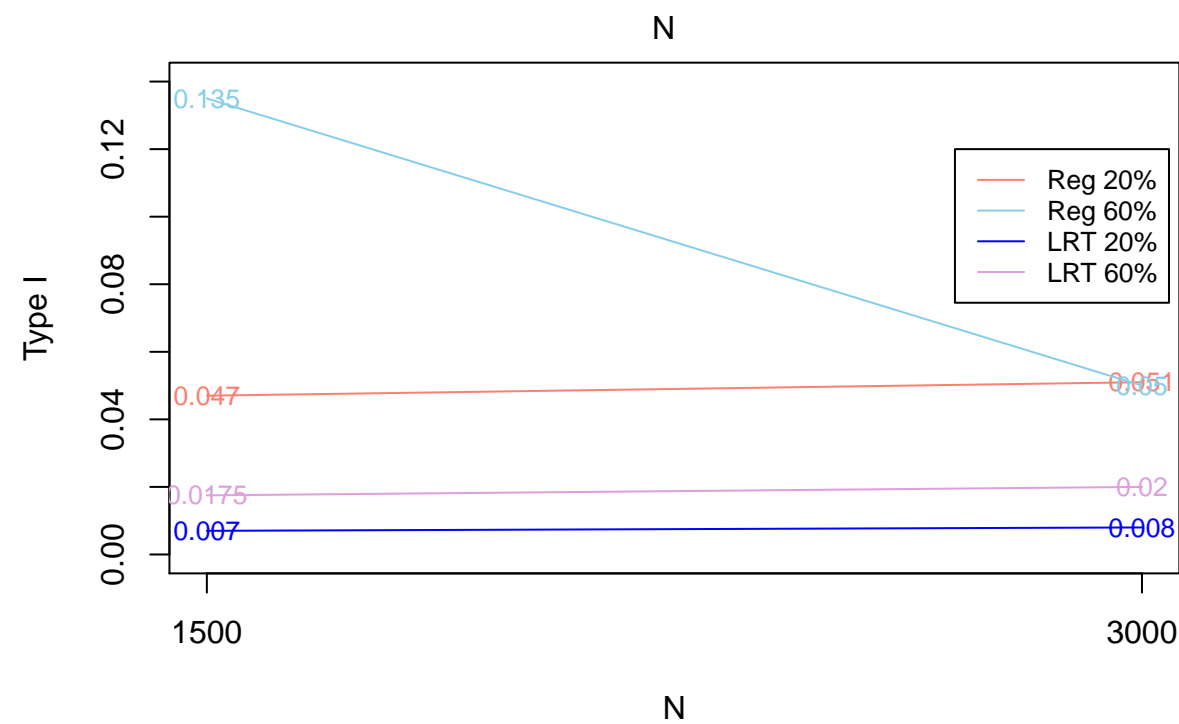
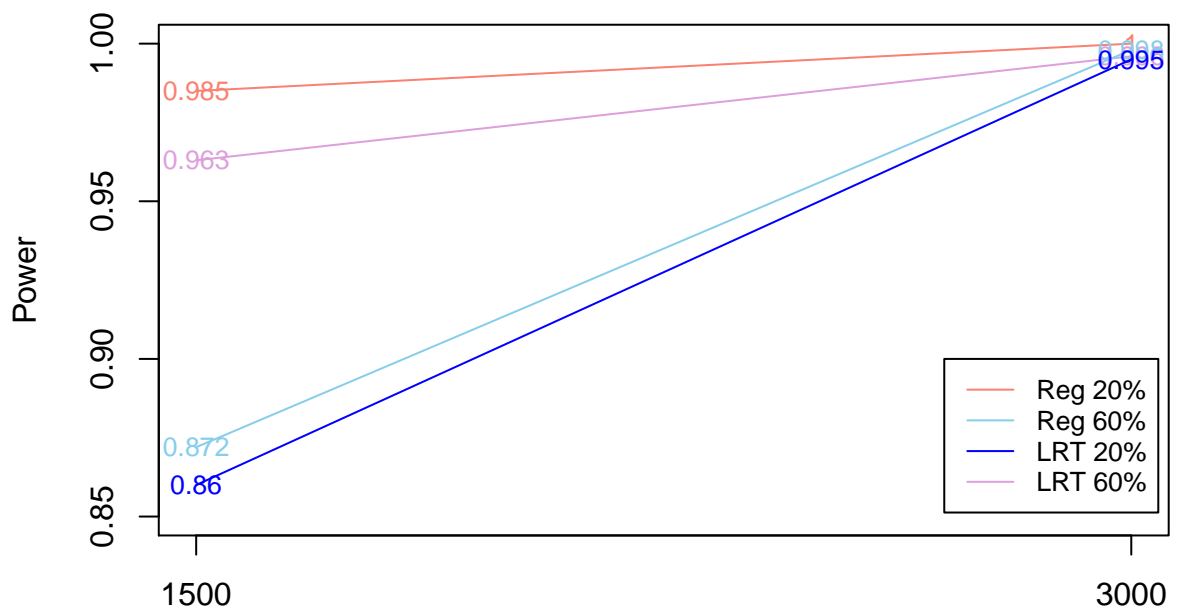
Group	Omnibus DIF
Power	0.941
Type I	0.0459

Table 34. Type I error and Power of mirt LRT

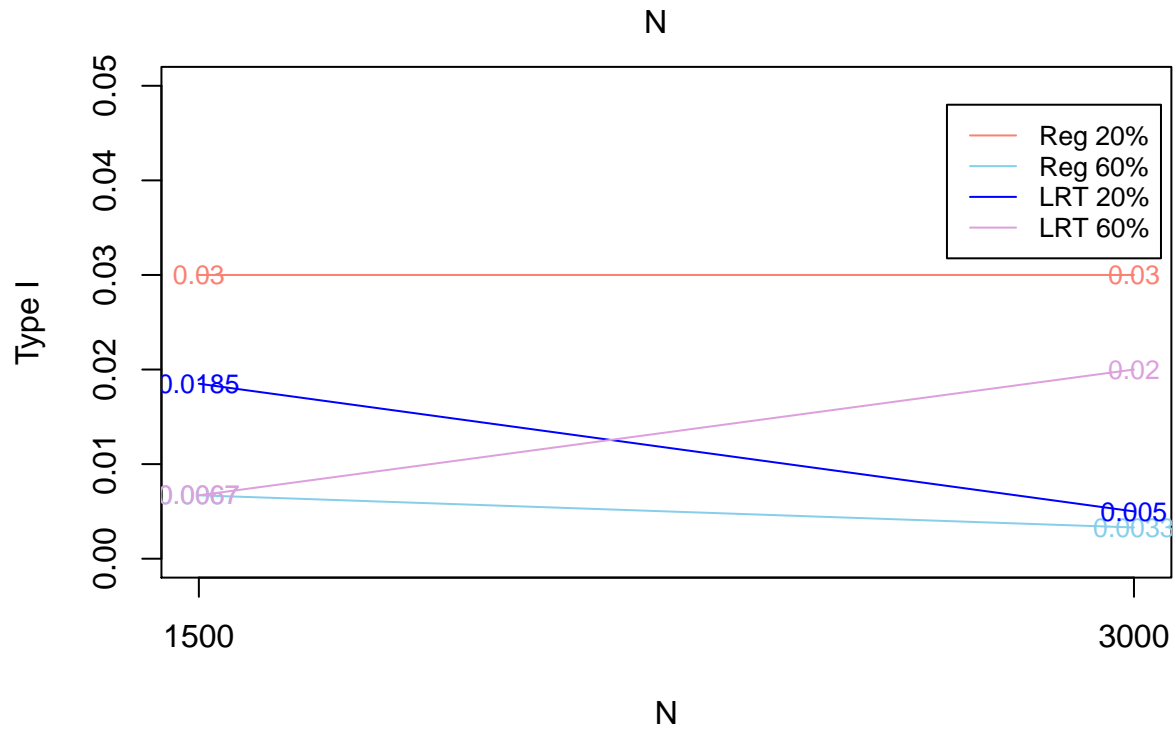
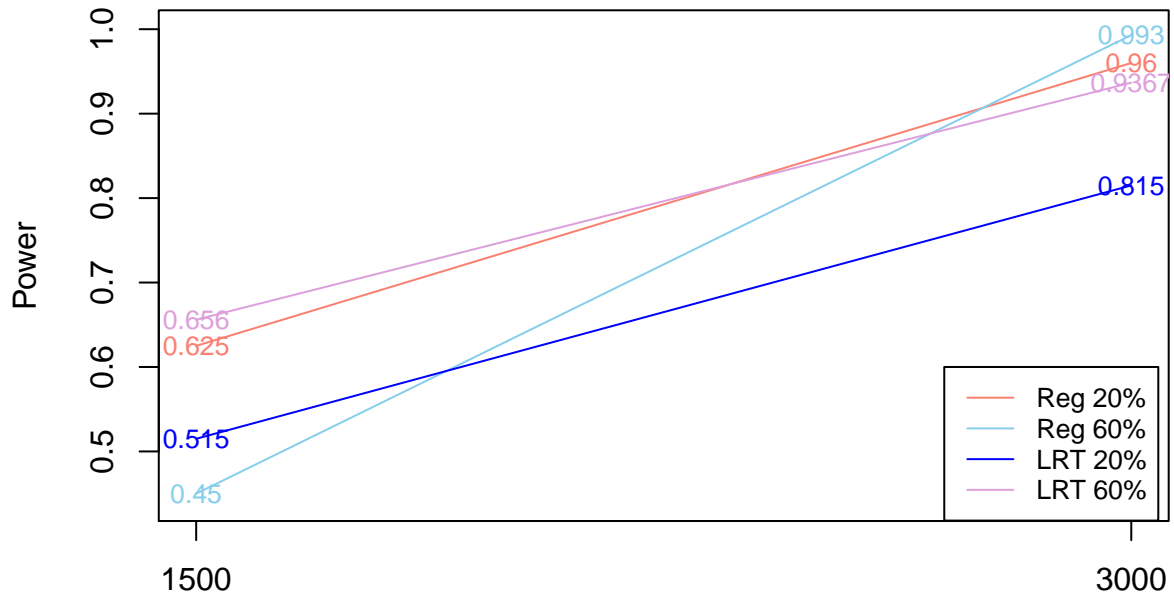
Group	Omnibus DIF
Power	0.985
Type I	0.023

## 8 Plots of Results

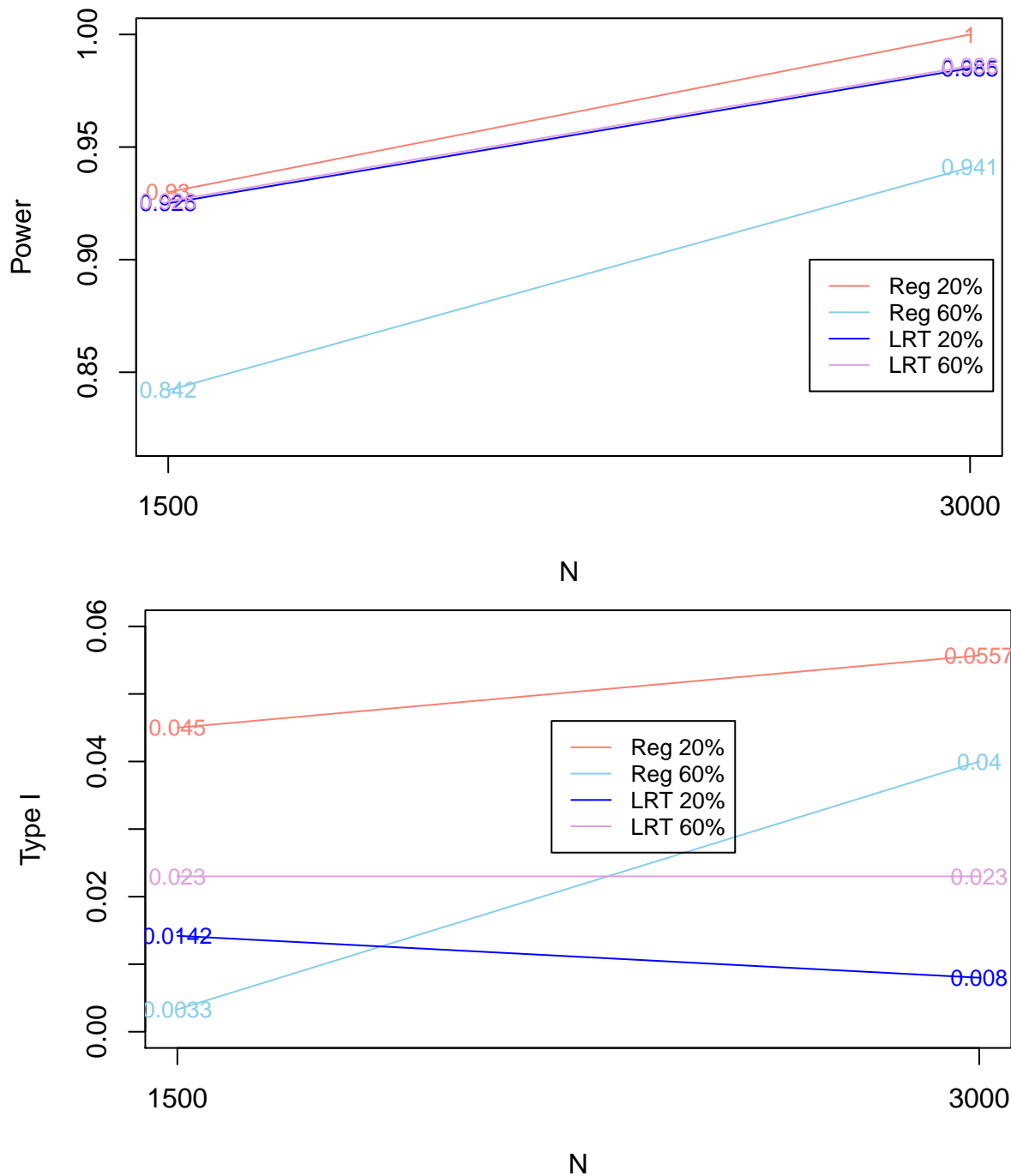
### 8.1 DIF only on intercept



## 8.2 DIF only on slope



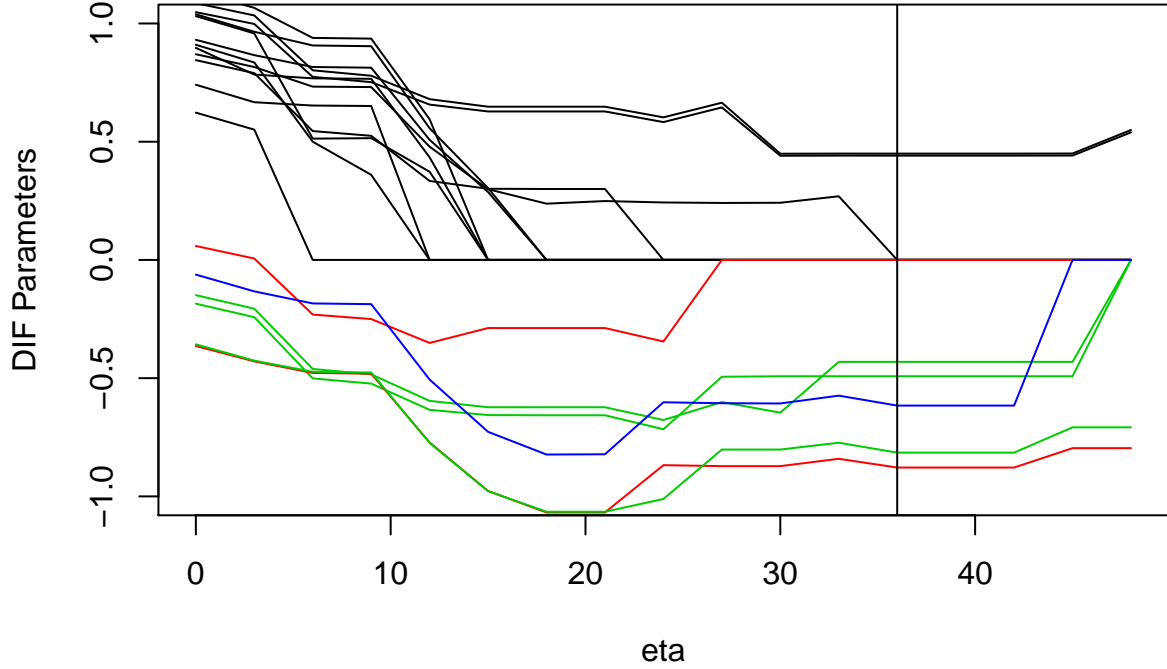
### 8.3 DIF on slope and intercept



## 9 High DIF Proportion Issue

When DIF proportion is higher than 50%, we found either our original LASSO or group LASSO algorithm work. Both of two methods would give us “reverse” type solutions– the true DIF parameters would be estimated as zero and non-DIF parameters would be non-zero. Below is a plot of the solution path of LASSO algorithm in 60% DIF on intercept condition. At the selected tuning parameter ( $\eta = 36$ ), most DIF items

(black lines) are shrunk to zero and non-DIF items (colored lines) are non-zero.



To solve this issue, we consider following methods:

1. “EMM” algorithm

From the above solution path plot, we can see when the tuning parameter is zero, the parameter estimations are good (DIF parameters are large and non-DIF parameters are close to zero). But the bias of LASSO estimates increases as tuning parameter increasing. So we considered adding a re-estimation step, that is, perform another M-step without penalty term after each EM cycle. The results of simulation 2, 4, 6, 8, 10 and 12 in this document are estimated by EMM algorithm and they look acceptable.

2. Using a different constraint

In section 4, we mentioned using anchor items for model identification. Those anchor items have no DIF on all focal groups. Here we can use a different constraint by setting “group-specific” anchors. Suppose we have two focal groups. Two anchor items can be used and each of them has DIF on one focal group but has no DIF on the other focal group.

To be more specific, recall

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

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

If the first two items are anchors, we want to set  $\beta_{j11} = 0$  and  $\beta_{j22} = 0$ , and  $\beta_{j21}$  and  $\beta_{j12}$  be non-zero and their magnitudes can be freely estimated. This constraint strongly relies on a priori knowledge of the items and the entire study but it works well when the DIF proportion is high.

3. Adaptive LASSO

Pick a  $\lambda > 0$ , and define the weight vector

$$\hat{w} = 1/|\hat{\beta}|^\lambda$$

( $\hat{\beta}$  can be the MLE).

$$(\hat{\mathbf{A}}, \hat{\mathbf{D}}, \hat{\beta}) = \operatorname{argmin}\{-\log M + \eta_n \sum_j^m \hat{w}_j \|\beta_j\|_1\} \quad (63)$$

Works.

4. SCAD penalty (Fan & Li, 2001)