# Improving Robustness in Q-Matrix Validation Using an Iterative and Dynamic Procedure

4 authors:

Pablo Nájera
Universidad Pontificia Comillas
**17** PUBLICATIONS   **93** CITATIONS

Miguel A. Sorrel
Universidad Autónoma de Madrid
**55** PUBLICATIONS   **705** CITATIONS

Jimmy de la Torre
The University of Hong Kong
**105** PUBLICATIONS   **4,675** CITATIONS

Francisco J Abad
Universidad Autónoma de Madrid
**141** PUBLICATIONS   **4,073** CITATIONS

Some of the authors of this publication are also working on these related projects:

Digital Literacy Assessment and Digital Citizenship View project

Q-matrix empirical validation procedures View project

5

6   Improving Robustness in Q-Matrix Validation using an Iterative and Dynamic Procedure

7

8   Pablo Nájera[a], Miguel A. Sorrel[a], Jimmy de la Torre[b], & Francisco José Abad[a]

9       [a]: Autonomous University of Madrid; [b]: The University of Hong Kong

10

11

12                                          Author Note

13      Pablo Nájera, Miguel A. Sorrel, and Francisco José Abad, Department of Social

14  Psychology and Methodology, Autonomous University of Madrid, Spain. Jimmy de la Torre,

15  Faculty of Education, The University of Hong Kong, Hong Kong.

20      Correspondence concerning this article should be addressed to Miguel A. Sorrel,

21  Department of Social Psychology and Methodology, Autonomous University of Madrid,

22  Ciudad Universitaria de Cantoblanco, Madrid 28049, Spain, e-mail: miguel.sorrel@uam.es.

23

24

25    **Improving robustness in Q-matrix validation using an iterative and dynamic procedure**

26                                         Abstract

27    In the context of cognitive diagnosis models, a Q-matrix reflects the correspondence between

28    attributes and items. The Q-matrix construction process is typically subjective in nature,

29    which may lead to misspecifications. All this can negatively affect the attribute classification

30    accuracy. In response, several methods of empirical Q-matrix validation have been developed.

31    The general discrimination index (GDI) method has some relevant advantages, such as the

32    possibility of being applied to several CDMs. However, the estimation of the GDI relies on

33    the estimation of the latent groups sizes and success probabilities, which is made with the

34    original (possibly misspecified) Q-matrix. This can be a problem, especially in those

35    situations in which there is a great uncertainty about the Q-matrix specification. To address

36    this, the present study investigates the iterative application of the GDI method where only one

37    item is modified at each step of the iterative procedure, and the required cutoff is updated

38    considering the new parameter estimates. A simulation study was conducted to test the

39    performance of the new procedure. Results showed that the performance of the GDI method

40    improved when the application was iterative at the item level and an appropriate cutoff point

41    was used. This was most noticeable when the original Q-matrix misspecification rate was

42    high, where the proposed procedure performed better 96.5% of the times. The results are

43    illustrated using Tatsuoka's fraction-subtraction dataset.

44    *Key words*: CDM, G-DINA, Q-matrix, validation, GDI.

45

46    **Improving robustness in Q-matrix validation using an iterative and dynamic procedure**

47    In the context of cognitive diagnosis assessment, cognitive diagnosis models (CDMs)

48    are latent class multidimensional statistical models that classify examinees as masters or non-

49    masters of different skills. Those skills are often referred to as *attributes*. Several CDMs have

50    been developed in the last years, which can be categorized as either reduced or general

51    models. The reduced models are the most specific ones; they provide low generalization but

52    high parsimony. The *deterministic input noise* and *gate* (DINA; Haertel, 1984; Junker &

53    Sijtsima, 2001), the *deterministic input noise* or *gate* (DINO; Templin & Henson, 2006), and

54    the *noisy input, deterministic output* and *gate* (NIDA; Maris, 1999; Junker & Sijtsima, 2001)

55    are some of the most widely known reduced models. Reduced models are usually preferred

56    because of the less number of parameter estimates and ease of interpretation. However, they

57    make strong assumptions about the data and model fit is therefore compromised. Reduced

58    models are nested in the general models, which allow for greater flexibility, but with more

59    demanding requirements (e.g., larger sample sizes). The *general diagnosis model* (GDM; von

60    Davier, 2005) and the *generalized DINA model* (G-DINA; de la Torre, 2011) are two

61    examples of general models. These models are preferred when there is not enough evidence to

62    assume a specific response process underlying the item responses.

63    The estimation of a CDM typically requires two inputs: the item responses of the

64    examinees and a Q-matrix (Tatsuoka, 1983). The Q-matrix is a $J$ (number of items) $\times K$

65    (number of attributes) matrix that reflects which attributes are measured by each item. Thus,

66    each item will have a q-vector ($\mathbf{q}_i$), in which each q-entry ($q_{ik}$) will adopt a value of 1 or 0

67    denoting if attribute $k$ is relevant for correctly answering item $j$ or not, respectively.

68    The original Q-matrix construction process should have a theoretical foundation, and

69    thus it is usually performed after a literature review, by analyzing examinees' reports, or by

70    domain experts. These processes are subjective in nature and can lead to some

71    misspecifications in the Q-matrix. These Q-matrix misspecifications negatively affect the

72    estimation of the model parameters and the accuracy of the attribute profile classification

73    (Gao, Miller & Liu, 2017; Rupp & Templin, 2008). For this reason, in the last years, several

74    empirically-based methods of Q-matrix validation have been developed with the aim of

75    detecting and correcting misspecified entries in a Q-matrix.

76         The present paper will focus on the *general discrimination index* (GDI) method, also

77    known as the general method of Q-matrix validation, developed for the G-DINA framework

78    by de la Torre and Chiu (2016). The structure of the paper will be the following. First, the G-

79    DINA model will be briefly introduced, followed by a description of the GDI method and its

80    advantages and limitations. Second, an item-level iterative procedure for the GDI method is

81    proposed and described. Third, the performance of the iterative procedure is compared to that

82    of the GDI method by means of Monte Carlo simulation. Fourth, a real data illustration is

83    conducted. Finally, a discussion of the results is provided, as well as future research insights

84    and comments on the advantages and limitations of the proposed procedure.

85    **Review of the G-DINA model**

86         The G-DINA model (de la Torre, 2011) is a general, saturated CDM that subsumes

87    most of the reduced models (e.g., DINA, DINO, *A*-CDM). In its original formulation, the

88    probability of success can be decomposed into the sum of the effects due to the presence of

89    specific attributes and their interactions:

$$P\left(\boldsymbol{\alpha}_{lj}^{*}\right) = \delta_{j0} + \sum_{k=1}^{K_j^*} \delta_{jk}\alpha_{lk} + \sum_{k'=k+1}^{K_j^*}\sum_{k=1}^{K_j^*-1} \delta_{jkk'}\alpha_{lk}\alpha_{lk'} \ ... + \delta_{12...K_j^*}\prod_{k=1}^{K_j^*}\alpha_{lk}, \quad (1)$$

90    where $\boldsymbol{\alpha}_{lj}^{*}$ is the reduced attribute vector whose elements are relevant for solving the item $j$;

91    $\delta_{j0}$ is the intercept of item $j$; $\delta_{jk}$ is the main effect due to $\alpha_k$; $\delta_{jkk'}$ is the interaction effect due

92    to $\alpha_k$ and $\alpha_{k'}$; and $\delta_{12...K_j^*}$ is the interaction effect due to $\alpha_1, ... , \alpha_{K_j^*}$, where $K_j^*$ is the number

93    of attributes specified for item $j$.

**The GDI method of empirical Q-matrix validation**

The GDI method of empirical Q-matrix validation (de la Torre & Chiu, 2016) is a

generalization of the $\delta$-method (de la Torre, 2008) that was developed for the DINA model.

The GDI method has been shown to perform well under both reduced and general CDMs at

detecting and modifying misspecifications in the Q-matrix. Apart from its great flexibility and

generalization, this method is included in the GDINA package (Ma & de la Torre, 2018) of the

R software (R Core Team, 2018) with a low computational cost. This makes it one of the

most accessible and easily applicable methods.

This validation method relies on the general discrimination index (GDI; usually

represented as $\varsigma_j^2$), which is the variance of the probabilities of success of the different latent

groups that are possible for an item weighted by the posterior distribution of those groups:

$$\varsigma_j^2 = \sum_{l=1}^{2^{K_j^*}} \omega(\boldsymbol{\alpha}_{lj}^*)\left[P(\boldsymbol{\alpha}_{lj}^*) - \bar{P}(\boldsymbol{\alpha}_{lj}^*)\right]^2 \tag{2}$$

where $2^{K_j^*}$ is the number of possible latent groups for item $j$, $\omega(\boldsymbol{\alpha}_{lj}^*)$ is the posterior

probability of examinees in group $\boldsymbol{\alpha}_{lj}^*$, $P(\boldsymbol{\alpha}_{lj}^*)$ is the probability of success for examinees in

this group, and $\bar{P}(\boldsymbol{\alpha}_{lj}^*)$ is the weighted mean probability of success across all the $2^{K_j^*}$ possible

latent groups for item $j$.

The method is based on the rationale that the correctly specified q-vector will lead to

the highest possible item discrimination value; that is, the correct q-vector for an item will be

the one that maximizes $\varsigma_j^2$. When comparing nested q-vectors, the specification of more

attributes in the q-vector will lead to a higher $\varsigma_j^2$, and thus a criterion needs to be included so

that the suggested q-vector for all items is not the one containing all the attributes ($\varsigma_{\mathbf{q}_j^{1:K}}^2$). De

la Torre and Chiu (2016) defined the *proportion of variance accounted for* (PVAF), which is

computed as $\text{PVAF}_{jc} = \varsigma_{\boldsymbol{q}_j^c}^2 / \varsigma_{\mathbf{q}_j^{1:K}}^2$, where $c$ reflects each of the $2^{K^*} - 1$ possible q-vectors

116    (note that the zero q-vector, with no attributes specified, is not plausible). The inclusion of

117    spurious attributes is prevented by determining a cutoff point ($\epsilon$, also referred as *EPS* for

118    *epsilon*), so the *suggested* q-vector would be the simplest one (i.e., the one with less attributes

119    specified) among those that fulfill PVAF > *EPS*.

120    Despite the good performance of the validation method, the original study was not

121    without limitations, as de la Torre and Chiu (2016) noted. For instance, the authors did not

122    justify the election criterion for the value of the *EPS*, which was set to 0.95. This aspect of the

123    method was examined by Nájera, Sorrel, and Abad (2019), who found that the GDI method

124    showed a good performance under a wide set of conditions, given that an optimal *EPS* for

125    each specific condition was used. Specifically, they provided a predictive formula for the

126    optimal *EPS* as a function of the average item quality (*IQ*), the sample size (*N*), and the

127    number of items (*J*):

$$EPS = \text{inv.logit}(-0.405 + 2.867 \cdot IQ + 4.840 \cdot 10^{-4} \cdot N - 3.316 \cdot 10^{-3} \cdot J), \quad (3)$$

128    where *inv.logit*($\cdot$) represents the inverse function of the logit function, computed as

129    $\exp(x)/(1 + \exp(x))$. *IQ* is computed as the average item quality ($IQ = \frac{1}{J}\sum_{j=1}^{J} IQ_j$), where

130    *IQ_j* is the difference in the probability of success between the latent group that possesses all

131    the relevant attributes specified in item *j*, $P_j(\mathbf{1})$, and the one with none of them, $P_j(\mathbf{0})$.

132    There is another aspect of the GDI method that deserves specific attention. When

133    computing $\varsigma_j^2$, the method assumes that the Q-matrix is correctly specified: $\varsigma_j^2$ relies on the

134    estimation of the latent group sizes and their success probabilities, which are estimated using

135    the provisional (misspecified) Q-matrix. As the authors point out, "it would be difficult, if not

136    impossible, for the same experts to correctly specify all the entries of the Q-matrix,

137    particularly when the test is long. Consequently, (b) [this assumption] is expected to always

138    be violated" (de la Torre & Chiu, 2016, p. 258). The authors state that the violation of the

139     assumption "does not automatically invalidate the viability of the proposed method. […] the

140     proposed method appears to be robust when the misspecifications in the Q-matrix is

141     controlled at a reasonable rate, which justifies the usefulness of the method in practice" (de la

142     Torre & Chiu, 2016, p. 258). According to the favorable results found by them with 5% of

143     misspecifications, and with 10% of misspecifications by Nájera et al. (2019), the method

144     seems indeed to be robust when the misspecification rate is low.

145         However, relying on the experts to make few mistakes while specifying the Q-matrix

146     is another assumption that may not always be realistic or, at least, will remain uncertain. It is

147     reasonable to think that different knowledge domains may vary in terms of Q-matrix

148     specification difficulty. For instance, the Q-matrix of a scholastic exam of mathematical

149     operations seems easier to specify (e.g., "*8 + 3 × 2*", would be easily detected as measuring,

150     for example, "sum" and "multiplication", but not "subtraction" or "division") than the Q-

151     matrix of a reading comprehension test, a clinical diagnostic test, or a test assessing students'

152     competencies (e.g., Sorrel et al. [2016] reported lower inter-rater reliability for more abstract

153     attributes like *"Study attitudes*" compared to attributes easier to objectivize like "*Helping*

154     *others*"). In fact, the Q-matrix of the popular fraction subtraction data set (Tatsuoka, 1990),

155     which does not belong to a particularly ambiguous knowledge domain, is still controversial

156     (Kang, Yang, & Zeng, 2019). Thus, the degree of uncertainty involved in the process could

157     reasonably be higher than what has been assumed, especially when the response processes of

158     the knowledge domain are somehow subjectively defined. Some authors have taken this point

159     under consideration, and have used in their simulation studies misspecification rates up to

160     40% (e.g., Wang et al., 2018). In light of the above, it is expected that the GDI method

161     performance will be compromised if the misspecification rate is reasonably high, since the

162     noise entered by the large number of misspecified q-entries can disrupt the calculation of $\varsigma_j^2$.

163     **Iterative Q-matrix validation methods**

164        One way of mitigating the pernicious effects that the violation of the true Q-matrix

165        assumption may provoke is to apply the validation method with an iterative procedure. Some

166        validation methods follow this rationale. The *iterative modified sequential search algorithm*

167        (IMSSA; Terzi & de la Torre, 2018a) and the *iterative general discrimination index* method

168        (iGDI; Terzi, 2017; Terzi & de la Torre, 2018b) are two validation methods in which all

169        proposed q-vector modifications are introduced in the Q-matrix in each iteration. In this

170        sense, they can be referred to as *test-level* iterative methods. On the other hand, the *Q-matrix*

171        *refinement method* (QRM; Chiu, 2013) and the data-driven approach proposed by Liu, Xu,

172        and Ying (2012) update the Q-matrix after each q-vector modification; that is, they modify

173        only one item in each iteration. Thus, they can be referred to as an *item-level* iterative method.

174        Even though *test-level* iterative methods can improve the performance of non-iterative

175        methods, it may be more precise to apply the iterative procedure at the item level. At the test-

176        level iteration, the first step will introduce several modifications based on the original and

177        presumably misspecified Q-matrix, and thus the probability of introducing wrong

178        modifications will be high. At the item-level, only the first item will be modified based on the

179        information of the original Q-matrix, while the rest of the items will be modified based on

180        progressively better specified Q-matrices. In the context of the GDI method, this will result in

181        a better recovery of $\varsigma_j^2$ and a more precisely predicted *EPS* as the iterations take place.

182        In light of the above, an optimal method should take into consideration the following

183        desired characteristics: first, it should be conducted iteratively; second, the iterations should

184        be applied at the item level; third, if a cutoff point is required, it should be selected by

185        empirical means and updated within each iteration; fourth, it should be applicable to both

186        reduced and general models. Based on this, it is expected that an item-level iterative

187        procedure based on the GDI method, applied with an optimal *EPS* that gets updated after each

188    iteration, will lead to promising results. The steps of the iterative procedure algorithm

189    evaluated in this paper are the following:

190        **Step 1**: Estimate the CDM according to the item responses and the provisional Q-

191        matrix (**Q**).

192        **Step 2**: Select the *EPS* value.

193        **Step 3**: Compute all items' $\varsigma_j^2$ (and PVAF) for each possible q-vector specification and

194        define, for each item, the set of *appropriate q-vector(s)*, which fulfill(s) PVAF > *EPS*.

195        **Step 4**: Select, for each item, the simplest element(s) among all the *appropriate q-*

196        *vectors*.

197            **4.1**: If there is only one element, then it is defined as the *suggested q-vector*.

198            **4.2**: If there are more than one element, the one with the highest PVAF is defined

199            as the *suggested q-vector*.

200        **Step 5**: Define, for each item, $\text{PVAF}_j^0$ as the PVAF of the *provisional q-vector*

201        specified in **Q**, and $\text{PVAF}_j^*$ as the PVAF of the *suggested q-vector*.

202        **Step 6**: Calculate all items' $\Delta\text{PVAF}_j$, defined as $\Delta\text{PVAF}_j = \left|\text{PVAF}_j^* - \text{PVAF}_j^0\right|$.

203        **Step 7**: Define the *hit item* as the item with the highest $\Delta\text{PVAF}_j$.

204        **Step 8**: Update **Q** by changing the *provisional q-vector* by the *suggested q-vector* of

205        the *hit item*.

206        **Step 9**: Iterate over Steps 1 to 8 until $\sum_{j=1}^{J} \Delta\text{PVAF}_j = 0$.

207        Step 2 and Steps 6 and 7 are of special relevance for the iterative procedure. Step 2

208    dictates which q-vectors are going to become *appropriate q-vectors* in Step 3 and,

209    consequently, which q-vector is going to become the *suggested q-vector* in Step 4. If the *EPS*

210    value is improperly chosen, the *suggested q-vectors* will be more likely to be incorrect. Thus,

211    each iteration will probably increase the distance between the provisional Q-matrix and the

212    true Q-matrix in a sort of "snowball" effect (i.e., errors will lead to more errors), and the $\varsigma_j^2$

213    will be worse specified. Hence, it is very important that the *EPS* election criterion is not

214    arbitrary. The predictive formula provided by Nájera et al. (2019; see Equation 3) showed a

215    good performance under a wide range of conditions. Furthermore, it can be easily

216    implemented in the iterative procedure and entails an additional benefit: as the prediction

217    formula considers the average item quality (*IQ*), which is computed after the model is

218    estimated, the *EPS* in Step 2 can be updated after each iteration. Step 7 is also very important,

219    because the election of the *hit item* can be neither be at random. Especially in the first

220    iterations, in which the Q-matrix will presumably still have several misspecifications, the $\varsigma_j^2$

221    is going to be calculated with some error. Steps 6 and 7 are used to select, for each iteration,

222    the q-vector that is more likely to be misspecified. These steps should optimize the

223    performance of the iterative procedure by increasing the probabilities of properly modifying a

224    q-vector in each iteration. The iterations would stop when all the *provisional q-vectors* and

225    *suggested q-vectors* are equal.

226                                        **Simulation study**

227            A simulation study was conducted to test if the proposed iterative procedure for the

228    GDI method provides better results than the standard (non-iterative) procedure. Two

229    hypotheses were stated: a) the iterative procedure will show a better performance than the

230    standard procedure, especially when the misspecification rate is high, b) this will be true as

231    long as the *EPS* value is properly chosen, based on the predictive formula. The performance

232    of the iterative procedure based on an inappropriate *EPS* value is expected to be worse than

233    that of the standard procedure, due to the "snowball" effect previously described.

234    **Method**

235            *Design*. The examinees' responses were simulated under the G-DINA model. The

236    number of attributes was fixed at $K = 5$, and the underlying distribution of examinees'

237    attribute patterns was uniform. The number of examinees was fixed at $N = 1000$, the average

238    item quality at $IQ = 0.6$, and the number of items at $J = 30$. Those values are considered to

239    be medium levels of each factor in applied contexts (Nájera et al., 2019). Table 1 shows the

240    Q-matrix used to simulate the examinees' responses ($\mathbf{Q}_{\text{true}}$). The Q-matrix was used in the

241    paper of de la Torre and Chiu (2016). It contains the same number of one-, two- and three-

242    attribute items, and each attribute is measured by the same number of items. Its structure

243    satisfies the required conditions to be a complete (Köhn & Chiu, 2017, 2018) and identifiable

244    (Gu & Xu, in press a, in press b) Q-matrix. Three variables were studied: the proportion of

245    misspecified q-entries or misspecification rate ($MR$ = 0.1, 0.2, 0.3, 0.4), the application

246    procedure for the GDI method (iterative, standard), and the $EPS$ value (predicted $EPS$, 0.95).

247    Thus, a total of 16 conditions resulted after combining the different factor levels (4

248    misspecification rates $\times$ 2 GDI application procedures $\times$ 2 $EPS$ values).

249    Table 1
250    *Q-Matrix for the Simulated Data*

| Item | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ | $\alpha_5$ | Item | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ | $\alpha_5$ |
|------|-----------|-----------|-----------|-----------|-----------|------|-----------|-----------|-----------|-----------|-----------|
| 1 | 1 | 0 | 0 | 0 | 0 | 16 | 0 | 1 | 0 | 1 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 | 17 | 0 | 1 | 0 | 0 | 1 |
| 3 | 0 | 0 | 1 | 0 | 0 | 18 | 0 | 0 | 1 | 1 | 0 |
| 4 | 0 | 0 | 0 | 1 | 0 | 19 | 0 | 0 | 1 | 0 | 1 |
| 5 | 0 | 0 | 0 | 0 | 1 | 20 | 0 | 0 | 0 | 1 | 1 |
| 6 | 1 | 0 | 0 | 0 | 0 | 21 | 1 | 1 | 1 | 0 | 0 |
| 7 | 0 | 1 | 0 | 0 | 0 | 22 | 1 | 1 | 0 | 1 | 0 |
| 8 | 0 | 0 | 1 | 0 | 0 | 23 | 1 | 1 | 0 | 0 | 1 |
| 9 | 0 | 0 | 0 | 1 | 0 | 24 | 1 | 0 | 1 | 1 | 0 |
| 10 | 0 | 0 | 0 | 0 | 1 | 25 | 1 | 0 | 1 | 0 | 1 |
| 11 | 1 | 1 | 0 | 0 | 0 | 26 | 1 | 0 | 0 | 1 | 1 |
| 12 | 1 | 0 | 1 | 0 | 0 | 27 | 0 | 1 | 1 | 1 | 0 |
| 13 | 1 | 0 | 0 | 1 | 0 | 28 | 0 | 1 | 1 | 0 | 1 |
| 14 | 1 | 0 | 0 | 0 | 1 | 29 | 0 | 1 | 0 | 1 | 1 |
| 15 | 0 | 1 | 1 | 0 | 0 | 30 | 0 | 0 | 1 | 1 | 1 |

251        *Data generation*. The probabilities of success of the latent groups with all the relevant

252    attributes, $P_j(\mathbf{1})$, and the probabilities of success of the latent groups with none of them,

253    $P_j(\mathbf{0})$, were manipulated to generate the item's quality ($IQ_j$). Specifically, $P_j(\mathbf{1}) =$

254    $U(0.7, 0.9)$ and $P_j(\mathbf{0}) = U(0.1, 0.3)$, which results in average values of $\bar{P}(\mathbf{1}) \cong 0.8$ and

255    $\bar{P}(\mathbf{0}) \cong 0.2$, giving an average item quality of $IQ = \bar{P}(\mathbf{1}) - \bar{P}(\mathbf{0}) \cong 0.6$. For the other latent

256    groups (those with some of the relevant attributes), the probabilities of success were simulated

257    so that they increased as the number of mastered attributes grew (i.e., monotonicity

258     constraint). Thus, a latent group that masters more attributes than other will always have

259     higher probabilities of success.

260        Misspecifications in the Q-matrix were introduced randomly with two constraints:

261     first, all items measured at least one attribute, and second, the first five items were not

262     modified. This latter constraint ensured the completeness of the Q-matrix, by assuring that

263     each attribute had, at least, one single-attribute item measuring it (Köhn & Chiu, 2017, 2018).

264        A total of 200 data sets were generated for each of the conditions. For each data set,

265     the $IQ_j$ were generated according to the aforementioned uniform distribution, and a different

266     misspecified Q-matrix ($\mathbf{Q}_{miss}$) was produced. All simulations and CDM analyses were

267     performed in R software, using the `GDINA` package.

268        *Dependent variables*. Two different types of dependent variables were used to assess

269     the performance of the validation method. First, the Q-matrix recovery rate (QRR) was used

270     to measure the quality of the Q-matrix specification recovery. It reflects the number of q-

271     entries that the method correctly specifies divided by the total number of q-entries ($J \times K$).

272     Second, the proportion of correctly classified attributes (PCA) and the proportion of correctly

273     classified vectors (PCV) were used to reflect the accuracy of attribute profile classification

274     (Ma & de la Torre, 2018). The PCA measures the proportion of entries (i.e., attributes)

275     correctly classified in the $N \times K$ matrix of attribute profile classification, while the PCV

276     reflects the proportion of examinees' attribute profiles that are completely correctly classified

277     (i.e., correctly classified rows in the $N \times K$ matrix of attribute classifications). Please note that

278     the PCV is a stricter measure than the PCA, and will usually obtain lower values. These

279     accuracy measures are of high relevance, since they provide information about the impact of

280     the Q-matrix specification quality in the final output of a CDM.

281        When applying a Q-matrix validation method, the suggested Q-matrix might show

282     some attributes positions (i.e., columns) interchanged. The possibility of having interchanged

283 attributes increases as the misspecification rate is higher. Thus, for each replica, the suggested

284 Q-matrix was compared with $\mathbf{Q}_{true}$ by checking the similarity between both matrices'

285 columns. Specifically, the mean absolute difference between the columns was conducted, and

286 the suggested Q-matrix's attribute columns were presented in the order that minimized the

287 difference with the corresponding $\mathbf{Q}_{true}$ attribute columns. This process is akin to a domain

288 expert labelling the factors when interpreting a factor analysis, where the order of the factors

289 is arbitrary. In the present case, the domain expert will evaluate whether the attributes are

290 correctly labelled.

291 **Results**

292       Before describing the main results, a brief comment about the iterative process (when

293 using the predicted *EPS*) is provided. No convergence problems were registered during the

294 simulation study. Table 2 shows the average number of iterations and number of items

295 modified (with one or more modifications in their q-vector) for each misspecification rate

296 condition. As expected, both measures increased as the misspecification rate did. It is

297 important to note that the number of iterations is usually higher than the number of items

298 modified, given that one item can be modified several times during the iteration procedure.

299 One item can be more properly modified at a later moment of the procedure, when the rest of

300 the Q-matrix is better specified. On the other hand, information about the average *IQ* and *EPS*

301 is given in Table 3. As expected, the initial *IQ* (i.e., the one estimated with the misspecified

302 Q-matrix) rapidly decreased as the misspecification rate increased. However, after the

303 iterative procedure was completed, the final *IQ* was adequately recovered, even for the most

304 unfavorable condition (i.e., $MR = 0.4$). This had an impact on the predicted *EPS*, which also

305 showed an increase from the original misspecified Q-matrix to the final validated Q-matrix.

306       In the following results, the performance of the standard and iterative procedures, as

307 well as their interaction with the predicted *EPS* and the *EPS* of 0.95, will be described. Tables

308  4, 5, and 6 show the results for the different dependent variables and conditions of the

309  simulation study in conjunction with the results obtained with the true Q-matrix and the

310  misspecified Q-matrices, which serve as upper and lower baselines, respectively. The type of

311  misspecification error (under- or over-specification) is disaggregated in Table 4. Plots for the

312  distribution of the dependent variables across the 200 replicates per misspecification rate

313  condition are provided in the Online Appendix. The different tables presented here include the

314  median of the 200 replicates due to the existence of asymmetry in the results distributions.

315  Results regarding the QRR, the PCA, and the PCV were consistent and showed similar

316  patterns. Thus, unless otherwise indicated, results for the three measures are described

317  together.

318  Table 2
319  *Average Number of Iterations and of Modified Items*

|  | Number of iterations | | | | Number of items modified* | | | |
|---|---|---|---|---|---|---|---|---|
| MR | Mean | SD | Min | Max | Mean | SD | Min | Max |
| 0.1 | 16.9 | 2.4 | 10 | 24 | 14.6 | 2.1 | 9 | 20 |
| 0.2 | 23.2 | 2.8 | 17 | 31 | 19.4 | 1.9 | 14 | 23 |
| 0.3 | 29.4 | 5.0 | 20 | 53 | 22.4 | 1.8 | 18 | 27 |
| 0.4 | 35.3 | 5.3 | 26 | 62 | 24.2 | 1.6 | 19 | 28 |

320  *Note*. * = with one or more modifications in their q-vector. MR = misspecification rate. This
321  information refers to the iterative procedure in conjunction with the predicted *EPS*.

322  Table 3
323  *Average Item Quality (*IQ*) and Used* EPS

|  | IQ | | EPS | |
|---|---|---|---|---|
| MR | Initial | Final | Initial | Final |
| 0.1 | 0.545 | 0.574 | 0.824 | 0.836 |
| 0.2 | 0.481 | 0.567 | 0.795 | 0.833 |
| 0.3 | 0.421 | 0.549 | 0.765 | 0.825 |
| 0.4 | 0.369 | 0.531 | 0.738 | 0.817 |

324  *Note*. MR = misspecification rate. Initial *IQ* and *EPS* values are obtained with the original
325  misspecified Q-matrix. Final *IQ* and *EPS* values are obtained with the validated Q-matrix
326  after the iterative procedure (using the predicted *EPS*) is completed. Items were simulated
327  with an *IQ* of 0.60.

328  As can be seen from Tables 4 to 6, the iterative implementation used in conjunction

329  with the predicted *EPS* always led to the best results. The Q-matrix recovery was very close

330  to one when the initial misspecification rate was low (QRR = 0.940), and was still high even

331    when the initial misspecification rate was high (QRR = 0.893). This procedure achieved the

332    highest QRR among the four presented procedures in most of the replicates, especially as the

333    misspecification rate increased. Thus, the iterative-predicted *EPS* implementation obtained the

334    highest QRR 62% of the times (*MR* = 0.1), 85.5% (*MR* = 0.2), 93.5% (*MR* = 0.3), and 96.5%

335    (*MR* = 0.4). It is important to note that, in those replicas in which it did not obtained the

336    highest QRR, it still obtained a QRR close to the highest, with a maximum loss of 0.07

337    through all misspecification rates. On the other hand, it obtained a QRR up to 0.32 higher

338    than the next best procedure, which reflects the better overall Q-matrix recovery shown in

339    Table 4. According to the GDI method rationale, a higher *EPS* tends to suggest more complex

340    q-vectors (i.e., with more attributes specified), and vice versa; thus, in Table 4 it can be seen

341    that the *EPS* of 0.95 produced more over-specification errors, while the predicted *EPS*

342    produced more under-specifications. The accuracy measures obtained with the iterative-

343    predicted *EPS* procedure were generally close to the upper limit regardless of the

344    misspecification rate. This was especially true for PCA. The misspecification rate affected

345    more severely the rest of the procedures. For example, the range of the median PCA values

346    reported in Table 5 for the standard and iterative implementations used in conjunction with

347    the predicted *EPS* were 0.085 and 0.012, respectively.

348    Table 4
349    *Medians for the Q-Matrix Recovery Rate (QRR) Results*

| MR | $Q_{true}$ | $Q_{miss}$ | Predicted *EPS* | | *EPS* = 0.95 | |
|----|-----------|-----------|------|------|------|------|
| | | | std | ite | std | ite |
| 0.1 | 1 | 0.900 | **0.940** | **0.940** | 0.887 | 0.833 |
| | | (6, 9) | (8, 1) | (8, 0) | (1, 16) | (1, 24) |
| 0.2 | 1 | 0.800 | 0.907 | **0.933** | 0.827 | 0.780 |
| | | (13, 17) | (11, 3) | (9, 1) | (2, 24.5) | (1, 32.5) |
| 0.3 | 1 | 0.700 | 0.817 | **0.913** | 0.720 | 0.687 |
| | | (19, 26) | (17, 11) | (11, 2) | (6, 36) | (1, 46) |
| 0.4 | 1 | 0.600 | 0.740 | **0.893** | 0.627 | 0.610 |
| | | (26, 34) | (21, 18) | (13, 3) | (8.5, 47) | (0.5, 58) |

350    *Note*. MR = misspecification rate; $Q_{true}$ = true Q-matrix; $Q_{miss}$ = misspecified Q-matrix; std =
351    standard procedure; ite = iterative procedure. A grayscale has been used for interpretation
352    purposes. Highest QRRs among the validation methods for each MR are shown in bold.

Median values for the number of under- and over-specified q-entries, respectively, are shown in brackets. Q-matrices are formed by 150 q-entries.

Table 5

*Medians for the Proportion of Correctly Classified Attributes (PCA) Results*

| MR | $Q_{true}$ | $Q_{miss}$ | Predicted *EPS* | | *EPS* = 0.95 | |
|---|---|---|---|---|---|---|
| | | | std | ite | std | ite |
| 0.1 | 0.910 | 0.895 | **0.907** | **0.907** | 0.900 | 0.894 |
| 0.2 | 0.911 | 0.867 | 0.901 | **0.906** | 0.894 | 0.889 |
| 0.3 | 0.911 | 0.813 | 0.862 | **0.903** | 0.868 | 0.880 |
| 0.4 | 0.910 | 0.764 | 0.822 | **0.895** | 0.807 | 0.864 |

*Note.* MR = misspecification rate; $Q_{true}$ = true Q-matrix; $Q_{miss}$ = misspecified Q-matrix; std = standard procedure; ite = iterative procedure. A grayscale has been used for interpretation purposes. Highest PCAs among the validation methods for each MR are shown in bold.

Table 6

*Medians for the Proportion of Correctly Classified Vectors (PCV) Results*

| MR | $Q_{true}$ | $Q_{miss}$ | Predicted *EPS* | | *EPS* = 0.95 | |
|---|---|---|---|---|---|---|
| | | | std | ite | std | ite |
| 0.1 | 0.637 | 0.583 | **0.625** | **0.625** | 0.603 | 0.581 |
| 0.2 | 0.642 | 0.484 | 0.604 | **0.623** | 0.586 | 0.560 |
| 0.3 | 0.643 | 0.325 | 0.457 | **0.613** | 0.492 | 0.531 |
| 0.4 | 0.639 | 0.227 | 0.337 | **0.579** | 0.335 | 0.483 |

*Note.* MR = misspecification rate; $Q_{true}$ = true Q-matrix; $Q_{miss}$ = misspecified Q-matrix; std = standard procedure; ite = iterative procedure. A grayscale has been used for interpretation purposes. Highest PCVs among the validation methods for each MR are shown in bold.

The following comments can be made regarding the manipulated factors. First, as it was expected, for both application procedures (standard vs. iterative) and *EPS* values (predicted *EPS* vs. *EPS* = 0.95), results were worse as the misspecification rate increased. Second, for both the standard and iterative procedures, and in line with the conclusions of Nájera et al. (2019), the predicted *EPS* provided better results than the *EPS* of 0.95. Third, regarding the interaction between the application procedure and the *EPS* value, the iterative procedure showed a better performance than the standard procedure only when the predicted *EPS* was used. Results were very similar for both procedures when the misspecification rate was low (*MR* = 0.1), but, as the misspecification rate was higher, the differences between both procedures substantially increased favoring the iterative procedure. On the contrary, when the *EPS* of 0.95 was used, the QRR of the iterative procedure was lower for all misspecification rates. As previously stated, these results were expected, since an inappropriate *EPS* increases

377     the probability of selecting an incorrect suggested q-vector, enlarging the distance between

378     the provisional Q-matrix and the true Q-matrix, disrupting the calculation of $\varsigma_j^2$. However,

379     regarding the PCA and the PCV, the iterative procedure, in conjunction with the *EPS* of 0.95,

380     showed slightly worse results when the misspecification rate was low (*MR* = 0.1 or 0.2), but

381     outperformed the standard procedure when the misspecification rate was high (*MR* = 0.3 or

382     0.4). All this reflects the fact that both an iterative procedure and a dynamic optimal *EPS*

383     value are required in order to achieve optimal results.

384                  **Real Data Example**

385     **Data and Analysis**

386         In order to facilitate a direct comparison between the proposed procedure and the

387     original GDI method, we used the same dataset as de la Torre and Chiu (2016). It consists of

388     536 examinees' responses to 11 fraction-subtraction items (Tatsuoka, 1990) measuring four

389     attributes (see strategy *b* in Mislevy, 1996): (1) performing basic fraction-subtraction

390     operation, (2) simplifying/reducing, (3) separating whole number from fraction, and (4)

391     borrowing one from whole number to a fraction. Table 7 shows the initial Q-matrix for these

392     data, which is the same as the one used by de la Torre and Chiu (2016). A higher-order G-

393     DINA model (de la Torre & Douglas, 2004) was used to fit the data.

394     **Results**

395         Table 7 shows the Q-matrix suggested by the iterative procedure. Six q-entries

396     modifications were proposed, all of them switching from 1 to 0, and all of them involving

397     attribute 2, with the exception of attribute 1 in Item 1. These results are somewhat congruent

398     with those found by de la Torre and Chiu (2016), who reported three modifications in

399     attribute 2 (Items 4, 5, and 11). According to the results found in the simulation results, the

400     iterative procedure suggested a less complex Q-matrix (i.e., less attributes specified) than the

401     original GDI method (see Table 4).

402    Regarding the original Q-matrix, attribute 2 (simplifying/ reducing) seems to have

403 theoretical relevance to solve the modified items. However, it is important to note that it

404 shows a great collinearity with attributes 3 and 4; that is, almost every time attribute 2 is

405 required, attributes 3 and 4 are also required. The only time that attribute 2 appears without

406 attributes 3 or 4 is in Item 6, which is the only one that retains attribute 2 in the suggested Q-

407 matrix. Thus, even though this attribute makes theoretical sense and seems to be correctly

408 specified in the original Q-matrix, it cannot be properly separated from other attributes. Since

409 it cannot provide any additional value, it becomes an irrelevant attribute and almost

410 disappeared in the suggested Q-matrix.

411 Table 7
412 *Original and suggested Q-matrices for the fraction-subtraction data*

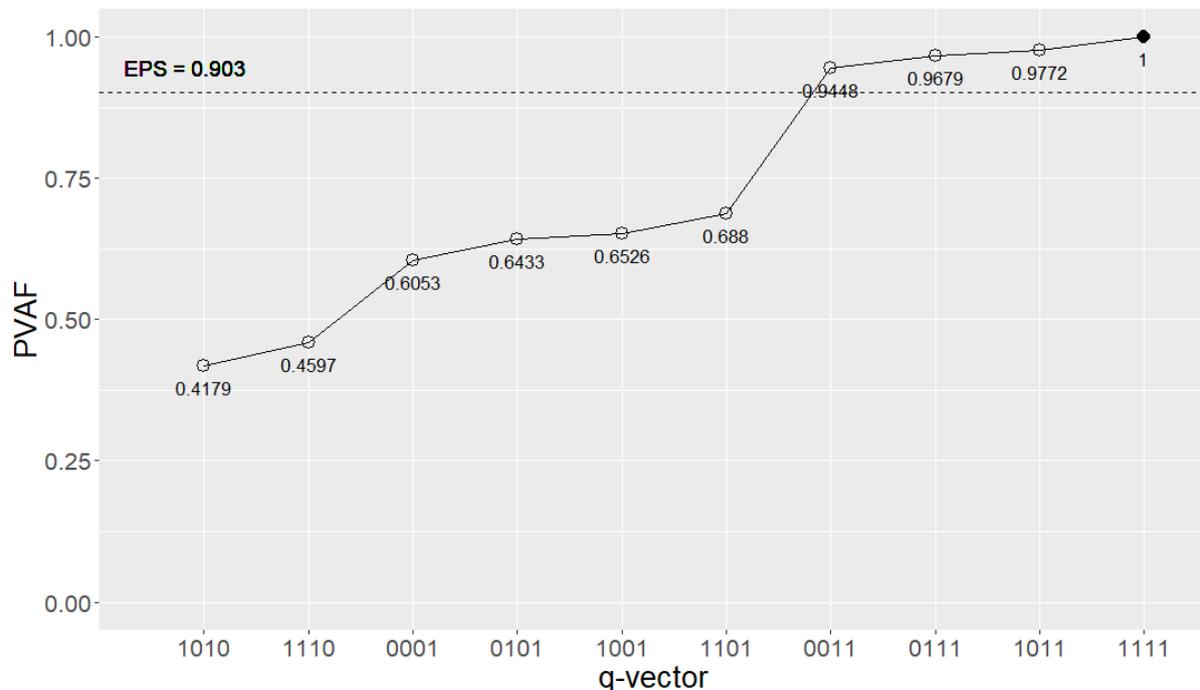| Item | | Original Q-matrix | | | | Suggested Q-matrix | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ | $\alpha_1$ | $\alpha_2$ | $\alpha_3$ | $\alpha_4$ |
| 1 | $3\frac{1}{2} - 2\frac{3}{2}$ | 1 | 1 | 1 | 1 | 0* | 0* | 1 | 1 |
| 2 | $\frac{6}{7} - \frac{4}{7}$ | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 3 | $3\frac{7}{8} - 2$ | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 4 | $4\frac{4}{12} - 2\frac{7}{12}$ | 1 | 1 | 1 | 1 | 1 | 0* | 1 | 1 |
| 5 | $4\frac{1}{3} - 2\frac{4}{3}$ | 1 | 1 | 1 | 1 | 1 | 0* | 1 | 1 |
| 6 | $\frac{11}{8} - \frac{1}{8}$ | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| 7 | $3\frac{4}{5} - 3\frac{2}{5}$ | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 8 | $4\frac{5}{7} - 1\frac{4}{7}$ | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 9 | $7\frac{3}{5} - \frac{4}{5}$ | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 |
| 10 | $4\frac{1}{10} - 2\frac{8}{10}$ | 1 | 1 | 1 | 1 | 1 | 0* | 1 | 1 |
| 11 | $4\frac{1}{3} - 1\frac{5}{3}$ | 1 | 1 | 1 | 1 | 1 | 0* | 1 | 1 |

413 *Note.* Q-entries modifications are highlighted with an asterisk.

414    Regarding Item 1, the first attribute is also removed in the suggested Q-matrix. This

415 item can be correctly solved by following different strategies:

416    (a) $3\frac{1}{2} - 2\frac{3}{2} = \frac{7}{2} - \frac{7}{2} = 0$ (attributes 1 and 4);

18

417 (b) $3\frac{1}{2} - 2\frac{3}{2} = 2\frac{3}{2} - 2\frac{3}{2} = 0$ (attributes 1, 3, and 4).

418 A *mesaplot* (Ma & de la Torre, 2018), which shows the PVAF related to each possible

419 q-vector specification, for Item 1 is presented in Figure 1. Four q-vectors (0011, 0111, 1011,

420 1111) clearly show a higher PVAF than the rest. Since their PVAF is higher than the *EPS*

421 (0.903), they form the set of appropriate q-vectors. The q-vector of 0011 is chosen as the

422 suggested q-vector because it is the simplest one. This attribute specification is related to

423 strategy (b), although attribute 1 is missing. A possible explanation to this could be that the

424 subtraction required in Item 1 may be a very easy operation that almost every examinee can

425 solve, since it involves two identical elements. As a consequence, attribute 1 would no longer

426 provide additional information. Nevertheless, these are modification suggestions, and domain

427 experts can seek among the appropriate q-vector in order to find the most suitable

428 specification. The last decision about the Q-matrix specification should rely on the judgment

429 of domain experts (de la Torre & Chiu, 2016).

430


431 *Figure 1*. Mesaplot for Item 1 of Tatsuoka's fraction-subtraction dataset included in Table 7.
432 The black dot represents the original q-vector specification (1111). The PVAF represents the
433 ratio of the GDI associated to a q-vector to the highest possible GDI that is obtained when all
434 the attributes are specified.

435          **Discussion**

436          CDMs rely on a correctly specified Q-matrix to provide an accurate classification of

437   examinees' attribute profiles. Domain experts are expected to specify the Q-matrix along with

438   a theoretical background, but they may commit some errors while doing so, especially when

439   the knowledge domain is particularly complex and ambiguous (e.g., mental pathologies,

440   reading comprehension, students' competencies). In this context, among the many Q-matrix

441   validation methods that have been developed in the last few years, de la Torre and Chiu

442   (2016) proposed the GDI method, which has some important advantages, such as its great

443   flexibility to be used with several reduced or general CDMs, its good performance at

444   modifying incorrectly specified q-vectors, and its low computational cost (Ma & de la Torre,

445   2018). Despite its benefits, the GDI method relies on the original Q-matrix, which may not be

446   correctly specified in most applied contexts. Although the method seemed robust to the

447   violation of this assumption when the Q-matrix misspecification rate was low, it is expected

448   to show a poorer performance when validating Q-matrices with more misspecifications.

449          The present paper evaluated an item-level iterative with dynamic *EPS* implementation

450   for the GDI method (this approach can be referred to as "ILD-GDI"). Considering past

451   research (e.g., Chiu, 2013; Liu et al., 2012; Nájera et al., 2019; Terzi & de la Torre, 2018ab),

452   we hypothesized that this implementation would lead to better results compared to the

453   existing procedures, especially when the misspecification rate is high. A simulation study was

454   conducted to test this hypothesis. Results showed that the new implementation did provide

455   better results. The gain obtained increased as the misspecification rate was higher.

456          The iterative procedure was hypothesized to have a poorer performance than the

457   standard procedure when used in conjunction with an inappropriate *EPS*. However, even

458   though the iterative-0.95 *EPS* (*ite95*) obtained a lower QRR than the standard-0.95 *EPS*

459   (*std95*), it provided better attribute profile classification results when the misspecification rate

20

460   was high (*MR* = 0.3 or 0.4). A tentative explanation of this result could be related to the type

461   of misspecification error. Some prior studies in the field (e.g., Gao, Miller, & Liu, 2017; Choi,

462   Templin, Cohen, & Atwood as cited in Kunina-Habenicht, Rupp, & Wilhelm, 2012) have

463   found that under-specifications have a greater impact in attribute profiles classification than

464   over-specifications. This effect is logically expected, since removing a parameter with a

465   substantive effect from a model might dramatically disrupt the probabilities of success of the

466   affected item; on the other hand, a spurious parameter added to the model may obtain a

467   marginal effect estimate, mitigating its impact (as long as the sample size is big enough to

468   produce stable parameter estimates).

469        This effect can explain the aforementioned results regarding *ite95* and *std95*. Table 4

470   shows the information regarding the Q-matrix recovery, disaggregated by specification error

471   type. On one hand, when *MR* = 0.1 or 0.2, *std95*'s QRR was higher than *ite95*'s. *Std95*'s PCA

472   and PCV were also higher than *ite95*'s. However, PCA differences were not as big as QRR

473   differences, since the higher amount of misspecifications in *ite95* were mainly over-

474   specifications, and both procedures had a similar number of under-specifications. On the other

475   hand, when *MR* = 0.3 or 0.4, *std95*'s QRR was still higher than *ite95*'s. However, *ite95*'s PCA

476   and PCV were higher than *std95*'s. Here, the QRR differences between both procedures were

477   smaller than those obtained with *MR* = 0.1 or 0.2. In addition, the higher amount of

478   misspecifications in *ite95* were mainly over-specifications, while *std95* obtained more under-

479   specifications. As previously stated, the latter might provoke a bigger disruption in the

480   posterior probabilities estimates, causing a worse attribute classification.

481        The explanation given above is certainly conditioned by the total number of

482   misspecifications. Under-specifications may have a bigger impact than over-specifications as

483   long as the total number of misspecifications remains at a similar range. The validation

484   procedure proposed in the present work (iterative in conjunction with the predicted EPS)

485     showed a higher number of under-specifications than *std95* and *ite95*; however, it showed a

486     much better performance in terms of Q-matrix specification recovery, which resulted in a

487     higher classification accuracy. It is important to note that other factors may have a relevant

488     role in modulating the relation between Q-matrix specification and attribute classification,

489     such as the number of different q-vectors represented in the Q-matrix (Rupp & Templin,

490     2008) and the identifiability of the Q-matrix (Gu & Xu, in press a, in press b).

491         Finally, a reviewer proposed examining whether the proposed procedure performs also

492     well when the underlying attribute's distribution is non-uniform. The performance of the

493     procedures under a multivariate normal distribution ($\rho = 0.25$; see Xu & Shang, 2018) and a

494     higher-order distribution ($\lambda_0 = (-1, -0.5, 0, 0.5, 1)$, $\lambda_{1k} = 1.5$; see de la Torre & Chiu, 2016)

495     are provided in the Online Appendix. It was observed that the pattern of results was very

496     similar to the ones obtained with the uniform distribution. Thus, the interpretation of the

497     findings do not differ according to the underlying attribute distribution, and the proposed

498     procedure still showed the best Q-matrix recovery and classification accuracy.

499         In conclusion, the ILD-GDI method proposed in this paper outperformed the original

500     method developed by de la Torre and Chiu (2016), as well as the method with the optimized

501     *EPS* value election (Nájera et al., 2019). The proposed procedure showed good performance

502     at detecting and modifying the Q-matrix even with a high misspecification rate (QRR $\geq$

503     0.893) and also at classifying attribute profiles (PCA $\geq$ 0.895; PCA$_{\mathbf{Q}\text{true}} \approx 0.910$), being the

504     only procedure that achieved a PCV higher than 0.5 under the worse misspecification rate

505     scenario (PCV $\geq$ 0.579; PCV$_{\mathbf{Q}\text{true}} \approx 0.640$). The iterative procedure's computation time was

506     short. On a laptop computer with four 2.2-GHz processors and 7 GB of RAM memory, the

507     average replica computation time under the worst condition (*MR* = 0.4) was 111 seconds.

508         The performance of the ILD-GDI method was also illustrated with Tatsuoka's

509     fraction-subtraction data. De la Torre and Chiu (2016) found that the standard GDI method

510    with an *EPS* of 0.95 proposed three modifications. These modifications were congruent with

511    the ones suggested by the ILD-GDI method. The suggestions of the ILD-GDI should be

512    considered rather than the GDI method's ones, since it provides a better recovery of the Q-

513    matrix, as shown in the simulation study. However, two consideration should be noticed.

514    First, even though Q-matrix validation methods are helpful in the search for the best possible

515    specified Q-matrix, some misspecifications may remain after their application. Second,

516    attribute positions in the Q-matrix are arbitrary just as factors are in a factor analysis; thus,

517    when two attributes (i.e., Q-matrix columns) have a similar specification through the items

518    and / or the number of misspecifications in the original Q-matrix is high, there exists the

519    possibility that the suggested Q-matrix shows interchanged positions for these attributes with

520    respect to the original Q-matrix. These considerations emphasize the role of domain experts in

521    the review of the validated Q-matrix. They should reject those suggested modifications that

522    lack a theoretical interpretation and check that the attributes maintain their original meaning.

523    Also, if they consider that several strategies can be followed to answer the items, multiple-

524    strategy models may be of help (e.g., de la Torre & Douglas, 2008; Ma & Guo, 2019). These

525    considerations may provide the most useful Q-matrix specification, since a tradeoff between

526    theoretical interpretation and data fit can be more easily achieved.

527         Further research is needed to extend the applicability of the ILD-GDI method. Even

528    though the performance of the GDI method was deeply studied under a wide range of

529    conditions by Nájera et al. (2019), the performance of the ILD-GDI method has only been

530    tested under a limited set of conditions. Further research would help to know whether it is

531    robust when the conditions are less favorable (e.g., small sample size, short test length, low

532    item quality). In this sense, other factors can be added to the study design, such as the number

533    of attributes or the underlying CDM (e.g., DINA).

534       Furthermore, it would be interesting to study whether the inclusion of model fit indices

535   to the iterative procedure could improve its performance. For instance, Kang et al. (2019)

536   used the item-level version of the RMSEA, which provided good results under the DINA

537   model. For the general CDMs framework, the Akaike's information criterion (AIC; Akaike,

538   1974) and the Bayesian information criterion (BIC; Schwarzer, 1976), which have been

539   previously used as fit indices in CDMs (e.g., Chen, de la Torre, & Zhang, 2013), could be

540   good candidates at selecting the *suggested q-vector*. One important drawback of this approach

541   would be the dramatic computational cost increment, since one additional model should be

542   estimated for each q-vector for each *hit item*. In this vein, the Wald test for model comparison

543   has also been recently used for Q-matrix validation under the sequential G-DINA model (Ma

544   & de la Torre, 2019).

**Final remarks**

546       The authors want to emphasize that empirical validation methods *suggest*

547   modifications, and cannot derive a *true* Q-matrix in empirical settings. The suggested Q-

548   matrix represents a model with empirical support. The purpose of Q-matrix validation should

549   not be to replace experts from the Q-matrix specification process, but to "provide

550   supplemental information for improving model-data fit, and consequently, increasing the

551   validity of inference from cognitive diagnosis assessments" (de la Torre & Chiu, 2016, p.

552   268). Especially in those contexts in which there is a certain degree of uncertainty involving

553   the Q-matrix, modification suggestions may help to understand which cognitive processes are

554   involved in responding each item. Also, as has been shown in the real data illustration,

555   validation methods can help detecting problems regarding the structure of the Q-matrix (e.g.,

556   attributes collinearity). Thus, we recommend applying three steps during the Q-matrix

557   specification process. First, construct the original Q-matrix with the help of domain experts.

558   In this step, the Delphi methodology can be of great help, facilitating the debate and

subsequent agreement between the judges (see Sorrel et al., 2016). It is also useful to track the degree of uncertainty involved in each q-entry during the process. Second, apply an empirical Q-matrix validation method, in order to detect any possible misspecifications made in the first step. Third, gather again the panel of experts to debate the theoretical viability of the suggested modifications and the meaning of the attributes after the process is completed. The degree of uncertainty involving each q-entry recorded in the first step can be of help at this point; a q-entry in which all experts showed a total agreement should probably not be modified even though the validation method suggests the opposite. In conclusion, the authors are of the opinion that the theory should be the main guide in the Q-matrix specification process. Empirical validation methods' role should be to support the domain experts' judgements.

<div align="center">**References**</div>

Akaike, H. (1974). A new look at the statistical identification model. *IEEE Transactions on Automated Control*, *19*, 716–723.

Chen, J., de la Torre, J., & Zhang, Z. (2013). Relative and absolute fit evaluation in cognitive diagnosis modeling. *Journal of Educational Measurement*, *50*, 123–140.

Chiu, C.-Y. (2013). Statistical refinement of the Q-matrix in cognitive diagnosis. *Applied Psychological Measurement, 37*(8), 598–618.

de la Torre, J. (2008). An empirically based method of Q-matrix validation for the DINA model: Development and applications. *Journal of Educational Measurement, 45*, 343–362.

de la Torre, J. (2011). The generalized DINA model framework. *Psychometrika, 76*, 179–199.

de la Torre, J., & Chiu, C.-Y. (2016). A general method of empirical Q-matrix validation. *Psychometrika, 81*(2), 253–273.

583    de la Torre, J., & Douglas, J. A. (2004). Higher-order latent trait models for cognitive

584        diagnosis. *Psychometrika*, *69*(3), 333–353.

585    de la Torre, J., & Douglas, J. A. (2008). Model evaluation and multiple strategies in cognitive

586        diagnosis: an analysis of fraction subtraction data. *Psychometrika*, *73*(4), 595–624.

587    Gao, M., Miller, M. D., & Liu, R. (2017). The impact of Q-matrix misspecification and model

588        misuse on classification accuracy in the generalized DINA model. *Journal of*

589        *Measurement and Evaluation in Education and Psychology, 8*(4), 391–403.

590    Gu, Y., & Xu, G. (in press a). Partial identifiability of restricted latent classes models. *Annals*

591        *of Statistics*. Retrieved from arXiv:1803.04353.

592    Gu, Y., & Xu, G. (in press b). Sufficient and necessary conditions for the identifiability of the

593        Q-matrix. *Statistica Sinica*. Retrieved from arXiv:1810.03819.

594    Haertel, E. (1984). An application of latent class models to assessment data. *Applied*

595        *Psychological Measurement*, *8*(3), 333–346.

596    Junker, B., & Sijtsma, K. (2001). Cognitive assessment models with few assumptions, and

597        connections with nonparametric IRT. *Applied Psychological Measurement, 25*, 258–

598        272.

599    Kang, C., Yang, Y., & Zeng, P. (2019). Q-matrix refinement based on item fit statistic

600        RMSEA. *Applied Psychological Measurement*, *43*(7), 527–542.

601    Köhn, H.-F., & Chiu, C.-Y. (2017). A procedure for assessing the completeness of the Q-

602        matrices of cognitively diagnostic tests. *Psychometrika*, *82*(1), 112–132.

603    Köhn, H.-F., & Chiu, C.-Y. (2018). How to build a complete Q-matrix for a cognitively

604        diagnostic test. *Journal of Classification*, *35*, 273–299.

605    Kunina-Habenicht, O., Rupp, A., & Wilhelm, O. (2012). The impact of model

606        misspecification on parameter estimation and item-fit assessment in log-linear

607        diagnostic classification models. *Journal of Educational Measurement*, *49*(1), 59–81.

608  Liu, J., Xu, G., & Ying, Z. (2012). Data-driven learning of Q-matrix. *Applied Psychological*

609      *Measurement*, *36*(7), 548–564.

610  Ma, W., & de la Torre, J. (2018). GDINA: The generalized DINA model framework. R

611      Package Version 2.0.8. Retrieved from https://cran.r-project.org/package=GDINA

612  Ma, W., & de la Torre, J. (2019). An empirical Q-matrix validation method for the sequential

613      generalized DINA model. *British Journal of Mathematical and Statistical Psychology*.

614      https://doi.org/10.1111/bmsp.12156

615  Ma, W., & Guo, W. (2019). Cognitive diagnosis models for multiple strategies. *British*

616      *Journal of Mathematical and Statistical Psychology*, *72*, 370–392.

617  Maris, E. (1999). Estimating multiple classification latent class models. *Psychometrika, 64*(2),

618      187–212.

619  Mislevy, R. J. (1996). Test theory reconceived. *Journal of Educational Measurement*, *33*(4),

620      379–416.

621  Nájera, P., Sorrel, M. A., & Abad, F. J. (2019). Reconsidering cutoff points in the general

622      method of empirical Q-matrix validation. *Educational and Psychological*

623      *Measurement*, *79*(4), 727–753.

624  R Core Team (2018). R (Version 3.4) [Computer Software]. Vienna, Austria: R Foundation

625      for Statistical Computing.

626  Rupp, A., & Templin, J. (2008). The effects of Q-matrix misspecification on parameter

627      estimates and classification accuracy in the DINA model. *Educational and*

628      *Psychological Measurement, 68*(1), 78–96.

629  Schwarzer, G. (1976). Estimating the dimension of a model. *Annals of Statistics*, *6*, 461–464.

630  Sorrel, M. A., Olea, J., Abad, F. J., de la Torre, J., Aguado, D., & Lievens, F. (2016). Validity

631      and reliability of situational judgement test scores: A new approach based on cognitive

632      diagnosis models. *Organizational Research Methods*, *19*(3), 506-532.

633    Tatsuoka, K. K. (1983). Rule space: An approach for dealing with misconception based on

634         item response theory. *Journal of Education Statistic, 20*, 345–354.

635    Tatsuoka, K. K. (1990). Toward an integration of item-response theory and cognitive error

636         diagnosis. In N. Frederiksen, R. Glaser, A. Lesgold, & Safto, M. (Eds.), *Monitoring*

637         *skills and knowledge acquisition* (pp. 453-488). Hillsdale, NJ: Lawrence Erlbaum.

638    Templin, J., & Henson, R. (2006). Measurement of psychological disorders using cognitive

639         diagnosis models. *Psychological Methods, 11*(3), 287–305.

640    Terzi, R. (2017). *New Q-matrix validation procedures* (Unpublished doctoral dissertation).

641         Rutgers, The State University of New Jersey, New Jersey, USA.

642    Terzi, R., & de la Torre, J. (2018a). An iterative method for empirically-based Q-matrix

643         validation. *International Journal of Assessment Tools in Education*, *5*, 248–262.

644    Terzi, R., & de la Torre, J. (2018b, April). *Two general iterative Q-matrix validation*

645         *procedures*. Paper presented at the meeting of the National Council of Measurement in

646         Education, New York, NY.

647    von Davier, M. (2005). A general diagnostic model applied to language testing data.

648         *Educational Testing Service, Research Report, RR-05-16*.

649    Wang, W., Song, L., Ding, S., Meng, Y., Cao, C., & Jie, Y. (2018). An EM-based method for

650         Q-matrix validation. *Applied Psychological Measurement*, 1–14.

651    Xu, G., & Shang, Z. (2018). Identifying latent structures in restricted latent class models.

652         *Journal of the American Statistical Association*, *113*(523), 1284–1295.