

¹ Factor Loading Recovery for Smoothed Tetrachoric Correlation Matrices

² Justin D. Kracht¹

³ ¹ University of Minnesota

Abstract

Researchers commonly use tetrachoric correlation matrices in item factor analysis. Unfortunately, tetrachoric correlation matrices are often indefinite (i.e., having one or more negative eigenvalues). These indefinite correlation matrices are problematic because the corresponding population correlation matrices they estimate are definitionally positive semidefinite (PSD; i.e., having strictly non-negative eigenvalues). Therefore, when used in procedures such as factor analysis, indefinite tetrachoric correlation matrices may result in poor estimates of factor loadings. Matrix smoothing algorithms attempt to remedy this problem by finding a PSD correlation matrix that is close, in some sense, to a given indefinite correlation matrix. However, little research has been done on the effectiveness of matrix smoothing for recovering the population correlation matrix, or for recovering factor loadings when smoothed matrices were used in exploratory factor analysis. In the present simulation study, indefinite tetrachoric correlation matrices were calculated from simulated binary data sets. Three matrix smoothing algorithms—the Higham (2002), Bentler-Yuan (2011), and Knol-Berger algorithms (1991)—were applied to the indefinite tetrachoric correlation matrices. Factor analysis was then conducted on the smoothed and unsmoothed correlation matrices. The results show that smoothed matrices were slightly better estimates of their population counterparts compared to unsmoothed indefinite correlation matrices. However, using smoothed compared to unsmoothed indefinite correlation matrices for item factor analysis did not meaningfully improve factor loading recovery. Matrix smoothing should therefore be considered only as a tool to facilitate factor analysis of indefinite correlation matrices and not as a statistical remedy for the root causes of matrix indefiniteness.

Keywords: matrix smoothing, item factor analysis, factor loading recovery, indefinite

27 Word count: X

28 Factor Loading Recovery for Smoothed Tetrachoric Correlation Matrices

29 Tetrachoric correlation matrices (Olsson, 1979) are used to estimate the correlations

30 between the normally-distributed, continuous latent variables often assumed to underlie

31 observed binary data. Therefore, tetrachoric correlation matrices are often recommended for

32 use in item factor analysis (i.e., factor analyses with binary or polytomous data) because the

33 common linear factor model requires the assumption that outcomes are continuous (Wirth &

34 Edwards, 2007). Unfortunately, tetrachoric correlation matrices are frequently indefinite,

35 having one or more negative eigenvalues (Bock, Gibbons, & Muraki, 1988; Wothke, 1993).

36 Indefinite tetrachoric correlation matrices are most likely to occur when computed from data

37 sets with many items, relatively small sample sizes, and extreme item loadings and

38 thresholds (Lorenzo-Seva & Ferrando, 2020). These Indefinite tetrachoric correlation

39 matrices are problematic because proper correlation matrices are, by definition, positive

40 semi-definite (PSD; i.e., having all eigenvalues greater than or equal to zero; Wothke, 1993).

41 Although indefinite correlation matrices resemble proper correlation matrices in many

42 ways—they are symmetric, have unit diagonals, and all off-diagonal elements less than or

43 equal to one in absolute value—it is impossible to obtain an indefinite matrix of Pearson

44 correlations from complete data. Thus, indefinite correlation matrices are improper estimates

45 of their corresponding population correlation matrices in the sense that they are not

46 included in the set of possible population correlation matrices.

47 Some researchers have suggested resolving the problem of indefinite tetrachoric

48 correlation matrices by obtaining a PSD correlation matrix that can be reasonably

49 substituted for an indefinite tetrachoric correlation matrix (e.g., Devlin, Gnanadesikan, &

50 Kettenring, 1975; Dong, 1985). This approach is often referred to as matrix smoothing, and

51 many algorithms developed for this purpose (referred to as matrix smoothing algorithms, or

52 simply smoothing algorithms) have been proposed in the psychometric literature and

53 elsewhere (Bentler & Yuan, 2011; Devlin et al., 1975; Dong, 1985; Fushiki, 2009; Higham,

54 2002; Knol & Berger, 1991; Li, Li, & Qi, 2010; Lurie & Goldberg, 1998; Qi & Sun, 2006).
55 However, despite the frequent occurrence of indefinite tetrachoric correlation matrices in
56 psychometric research (Bock et al., 1988, p. 261), the variety of smoothing algorithms
57 available, and suggestions to use matrix smoothing algorithms as a remedy to indefinite
58 tetrachoric correlation matrices (Bentler & Yuan, 2011; Knol & Berger, 1991; Wothke, 1993),
59 scant research has been done on the effectiveness of matrix smoothing algorithms in the
60 context of item factor analysis of indefinite tetrachoric correlation matrices (Lorenzo-Seva &
61 Ferrando, 2020). In one of the only published comparisons of this kind, Knol and Berger
62 (1991) investigated the effects of using smoothed compared to unsmoothed correlation
63 matrices in factor analysis and found no large differences in factor loading recovery. However,
64 this comparison was not a main focus of their study and only compared a small number of
65 indefinite matrices (10 indefinite correlation matrices with 250 subjects and 15 items).

66 Additionally, few studies have compared the *relative* performance of matrix smoothing
67 algorithms in the context of factor analysis (Debelak & Tran, 2013, 2016). Debelak and Tran
68 (2013) conducted a simulation study to determine which of three matrix smoothing
69 algorithms—the Higham alternating-projections algorithm (APA; 2002), the Bentler-Yuan
70 algorithm (BY; 2011), and the Knol-Berger (KB; 1991) algorithm—most often recovered the
71 underlying dimensionality when applied to indefinite tetrachoric correlation matrices prior to
72 parallel analysis (Horn, 1965). Debelak and Tran simulated binary data using a
73 two-parameter logistic (2PL) item response theory (IRT; Birnbaum, 1968; de Ayala, 2013)
74 model for one- and two-factor models with varying factor correlations, item difficulties, item
75 discriminations, numbers of items, and numbers of subjects. Debelak and Tran then
76 computed tetrachoric correlation matrices for each simulated binary data set. If a tetrachoric
77 correlation matrix was indefinite, the three aforementioned smoothing algorithms were
78 applied (resulting in three smoothed correlation matrices in addition to the indefinite
79 tetrachoric matrix). Finally, Debelak and Tran conducted parallel analysis using each of the
80 four correlation matrices to obtain estimates of dimensionality. Debelak and Tran concluded

81 that “[the] application of smoothing algorithms generally improved correct identification of
82 dimensionality when the correlation between the latent dimensions was 0.0 or 0.4 in our
83 simulations” (Debelak & Tran, 2013, p. 74). With respect to the relative performance of the
84 Higham, Bentler-Yuan, and Knol-Berger smoothing algorithms in this context, Debelak and
85 Tran concluded that there were “minor differences in the performance of the three smoothing
86 algorithms used in [the] study. In data sets with a clear dimensional structure...the
87 algorithm of Bentler and Yuan (2011) performed best” (Debelak & Tran, 2013, p. 74).

88 Following on these results, Debelak and Tran (2016) extended their previous simulation
89 design to evaluate the relative and absolute effectiveness of matrix smoothing algorithms
90 when applied to indefinite polychoric correlation matrices of ordered, categorical (i.e.,
91 polytomous) data prior to conducting a parallel analysis. As in their previous study, Debelak
92 and Tran used the accuracy of the parallel analysis dimensionality estimates as their
93 evaluation criterion. In addition to extending their design to consider polytomous data,
94 Debelak and Tran (2016) also considered factor models with either one or three major
95 common factors and either zero or forty minor common factors. The minor common factors
96 represented the effects of model approximation error; that is, the degree of model misfit
97 inherent to mathematical models of natural phenomena in general, and psychological models
98 in particular (MacCallum & Tucker, 1991; MacCallum, Widaman, Preacher, & Hong, 2001;
99 Tucker, Koopman, & Linn, 1969). Debelak and Tran concluded that the analysis of
100 smoothed polychoric correlation matrices generally led to more accurate results than the
101 analysis of indefinite polychoric correlation matrices. Moreover, they found that “methods
102 based on the algorithms of Knol and Berger, Higham, and Bentler and Yuan showed a
103 comparable performance with regard to the accuracy to detect the number of underlying
104 major factors, with a slightly better performance of methods based on the Bentler and Yuan
105 algorithm” (Debelak & Tran, 2016, p. 15).

106 Both Debelak and Tran (2013) and Debelak and Tran (2016) concluded that the

107 Bentler-Yuan (2011) smoothing algorithm led to the most accurate results (in terms of
108 dimensionality recovery) when applied to indefinite tetrachoric or polychoric correlation
109 matrices. However, neither study attempted to explain why the Bentler-Yuan algorithm led
110 to better dimensionality recovery relative to the other smoothing methods they investigated.
111 One intriguing possibility is that the smoothed correlation matrices produced by the
112 Bentler-Yuan algorithm were better approximations of the population correlation matrices
113 than either the smoothed matrices produced by the Knol-Berger (1991) and Higham
114 algorithms (2002) or the original indefinite tetrachoric or polychoric correlation matrices. If
115 this is true, one might also expect that Bentler-Yuan smoothed tetrachoric correlation
116 matrices will also lead to more accurate factor loading estimates compared to the
117 alternatives.

118 The purpose of the present study was to address two questions related to these
119 hypotheses. First, are smoothed indefinite tetrachoric correlation matrices better estimates
120 of their corresponding population correlation matrices than the original indefinite tetrachoric
121 correlation matrices and, if so, which smoothing method produces the best estimates?
122 Second, do smoothed indefinite tetrachoric correlation matrices lead to better factor loading
123 estimates compared to the unsmoothed tetrachoric matrices when used in exploratory factor
124 analysis and, if so, which smoothing algorithm leads to the best factor loading estimates? To
125 answer these questions, I conducted a simulation study in which I generated 124,346
126 indefinite tetrachoric correlation matrices from a variety of realistic data scenarios. Before
127 describing the simulation design, I first introduce tetrachoric correlations, the three matrix
128 smoothing algorithms under investigation, the common factor model, and the three factor
129 analysis algorithms included in this study.

130 **Tetrachoric Correlations**

131 A tetrachoric correlation is an estimate of the linear association between two
132 continuous, normally-distributed latent variables, y_1^* and y_2^* , obtained using dichotomous,

₁₃₃ observed manifestations of those variables, y_1 and y_2 . The variables y_1^* and y_2^* are assumed
₁₃₄ to follow a bivariate normal distribution,

$$\begin{pmatrix} y_1^* \\ y_2^* \end{pmatrix} \sim \mathcal{N} \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & r \\ r & 1 \end{pmatrix} \right],$$

₁₃₅ where r is the true correlation between y_1^* and y_2^* that is estimated by the tetrachoric
₁₃₆ correlation, \hat{r} . To compute the tetrachoric correlation, a 2×2 contingency table is first
₁₃₇ created using y_1 and y_2 as described in Brown and Benedetti (1977). If any of the cell
₁₃₈ frequencies in the contingency table are zero, those elements are replaced with 0.5 and the
₁₃₉ other elements adjusted to leave the marginal sums unchanged (Brown & Benedetti, 1977).
₁₄₀ The proportions of correct responses for y_1 and y_2 are represented by the marginals p_1 and
₁₄₁ p_2 . The standard normal deviate thresholds, τ_1 and τ_2 , used to dichotomize y_1^* and y_2^* are
₁₄₂ then estimated using $1 - \Phi(\hat{\tau}_1) = p_1$ and $1 - \Phi(\hat{\tau}_2) = p_2$, and solving for $\hat{\tau}_1$ and $\hat{\tau}_2$. Here,
₁₄₃ $\Phi(z)$ denotes the standard normal cumulative distribution function (Divgi, 1979). Because y_1^*
₁₄₄ and y_2^* are assumed to follow a bivariate normal distribution with correlation r , the joint
₁₄₅ probability of $(y_1 > \hat{\tau}_1, y_2 > \hat{\tau}_2)$ can be written as:

$$L(\hat{\tau}_1, \hat{\tau}_2, r) = \frac{1}{2\pi\sqrt{1-r^2}} \int_{\hat{\tau}_2}^{\infty} \int_{\hat{\tau}_1}^{\infty} \exp \left(-\frac{y_1^{*2} + y_2^{*2} - 2ry_1^*y_2^*}{2(1-r^2)} \right) dy_1^* dy_2^*. \quad (1)$$

₁₄₆ An estimate of r can then be obtained by setting Equation (1) equal to p_{11} (the
₁₄₇ observed proportion of correct responses for both y_1 and y_2) and solving for r using an
₁₄₈ iterative procedure. In particular, the Newton-Raphson method can be used to obtain
₁₄₉ successive approximations of r given an initial estimate, \hat{r}_0 :

$$\hat{r}_{i+1} = \hat{r}_i - \frac{L(\hat{\tau}_1, \hat{\tau}_2, \hat{r}_i) - p_{11}}{L'(\hat{\tau}_1, \hat{\tau}_2, \hat{r}_i)}, \quad (2)$$

150 where $L'(\hat{\tau}_1, \hat{\tau}_2, \hat{r}_i)$ is the first derivative of $L(\hat{\tau}_1, \hat{\tau}_2, \hat{r}_i)$ (Divgi, 1979). Iteration continues
 151 until convergence is achieved (when $\hat{r}_{i+1} - \hat{r}_i < \delta$ for some small value of δ) or until some
 152 maximum number of iterations occur. For p dichotomous variables, the $p \times p$ symmetric
 153 matrix \mathbf{R}_{Tet} is called the tetrachoric correlation matrix. The \mathbf{R}_{Tet} matrix has a unit diagonal
 154 and has off-diagonal elements consisting of pairwise tetrachoric correlation coefficients
 155 \hat{r}_{jk} , $j, k \in \{1, \dots, p\}$. Just as the tetrachoric correlation \hat{r}_{jk} estimates r_{jk} , the tetrachoric
 156 correlation matrix \mathbf{R}_{Tet} estimates the $p \times p$ population correlation matrix, \mathbf{R}_{Pop} , which is
 157 symmetric with off-diagonal elements r_{jk} , and a unit diagonal.

158 **Matrix Smoothing Algorithms**

159 **Higham Alternating Projections Algorithm (APA; 2002).** The matrix
 160 smoothing algorithm proposed by Higham (2002) seeks to find the closest PSD correlation
 161 matrix to a given indefinite correlation matrix. In this context, closeness is defined as the
 162 generalized Euclidean distance (Banerjee & Roy, 2014, p. 492). Higham's algorithm (2002)
 163 uses a series of alternating projections to locate the PSD correlation matrix (\mathbf{R}_{APA}) closest
 164 to a given indefinite correlation matrix (\mathbf{R}_{-}) of the same order. The algorithm works by first
 165 projecting \mathbf{R}_{-} onto the set of symmetric, PSD $p \times p$ matrices, \mathcal{S} . The resulting candidate
 166 matrix is then projected onto the set of symmetric $p \times p$ matrices with unit diagonals, \mathcal{U} . The
 167 series of projections repeats until the algorithm converges to a matrix, \mathbf{R}_{APA} , that is PSD,
 168 symmetric, and has a unit diagonal, or until the maximum number of iterations is exceeded.

169 Specifically, Higham's algorithm (2002) consists of alternating projection functions, P_U ,
 170 the projection onto \mathcal{U} , and P_S , the projection onto \mathcal{S} . For some symmetric $\mathbf{A} \in \mathbb{R}^{p \times p}$ with
 171 elements a_{ij} ,

$$P_U(\mathbf{A}) = (p_{ij}), \quad p_{ij} = \begin{cases} a_{ij}, & i \neq j \\ 1, & i = j. \end{cases} \quad (3)$$

172 Stated simply, $P_U(\mathbf{A})$ replaces all elements of the diagonal of \mathbf{A} with ones. The projection
 173 onto \mathcal{S} is less straightforward. Higham (2002) outlines the steps as follows. For some

₁₇₄ symmetric, indefinite matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$, let $\mathbf{A} = \mathbf{V}\Lambda_{-}\mathbf{V}^T$ be the eigendecomposition of \mathbf{A} ,
₁₇₅ where \mathbf{V} is the orthonormal matrix of eigenvectors and $\Lambda_{-} = \text{diag}(\lambda_i)$ is a diagonal matrix
₁₇₆ with the eigenvalues of \mathbf{A} , $\lambda_i, i \in \{1, \dots, p\}$, ordered from largest to smallest on the diagonal
₁₇₇ ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p, \lambda_p < 0$). Also let $\Lambda_{+} = \text{diag}(\max(\lambda_i, 0))$. Then the projection of \mathbf{A}
₁₇₈ onto \mathcal{S} can be written as

$$P_S(\mathbf{A}) = \mathbf{V}\Lambda_{+}\mathbf{V}^T. \quad (4)$$

₁₇₉ Starting with $\mathbf{A} = \mathbf{R}_{-}$, \mathbf{R}_{APA} can be obtained by repeatedly applying the operation
₁₈₀ $\mathbf{A} \leftarrow P_U(P_S(\mathbf{A}))$ until convergence occurs or until some maximum number of iterations is
₁₈₁ reached (Higham, 2002, p. 337).

₁₈₂ **Bentler-Yuan Algorithm (BY; 2011).** The Bentler-Yuan (2011) smoothing
₁₈₃ algorithm is based on minimum-trace factor analysis (MTFA; Bentler, 1972; Jamshidian &
₁₈₄ Bentler, 1998). MTFA seeks to find optimal communality estimates such that unexplained
₁₈₅ common variance is minimized. This minimization is subject to two constraints. First, the
₁₈₆ diagonal matrix of unique variances is constrained to be positive semidefinite (PSD). Second,
₁₈₇ the matrix formed by replacing the diagonal elements of the observed covariance matrix with
₁₈₈ the estimated communalities is also constrained to be PSD. In contrast with the Higham
₁₈₉ algorithm (2002), the Bentler-Yuan algorithm does not seek to minimize some criterion.
₁₉₀ Instead, the algorithm uses MTFA to identify Heywood cases (i.e., communality estimates
₁₉₁ greater than or equal to one and, consequently, negative or zero uniqueness variance
₁₉₂ estimates; Dillon, Kumar, & Mulani, 1987). The Bentler-Yuan algorithm then rescales the
₁₉₃ rows and columns of \mathbf{R}_{-} corresponding to these Heywood cases to produce a smoothed, PSD
₁₉₄ correlation matrix, \mathbf{R}_{BY} . More specifically, the algorithm first conducts an MTFA using \mathbf{R}_{-} .
₁₉₅ Using the results of the MTFA, a diagonal matrix, \mathbf{H} is constructed containing the estimated
₁₉₆ communalities as diagonal elements. Next, another diagonal matrix, Δ^2 , is constructed with
₁₉₇ elements δ_i^2 where $\delta_i^2 = 1$ if $h_i < 1$ and $\delta_i^2 = k/h_i$ otherwise (where $k < 1$ is some constant).
₁₉₈ Finally, the smoothed, PSD correlation matrix $\mathbf{R}_{\text{BY}} = \Delta \mathbf{R}_0 \Delta + \mathbf{I}$ is obtained, where \mathbf{R}_0 is
₁₉₉ \mathbf{R}_{-} with diagonal elements replaced by zeroes and \mathbf{I} is an identity matrix that ensures that

200 \mathbf{R}_{BY} has a unit diagonal.

201 Similar to the Higham algorithm, the Bentler-Yuan algorithm sometimes fails to
202 produce a PSD correlation matrix. This can happen either when (a) the MTFA algorithm
203 fails to converge or (b) when k is too large and does not shrink the targeted elements of the
204 indefinite correlation matrix enough for the matrix to become PSD. To help with this
205 non-convergence, I used the modified Bentler-Yuan algorithm implementation provided by
206 the `smoothBY()` function in the R *fungible* package (Waller, 2019) to adaptively select an
207 appropriate k . The k parameter was initialized at $k = 0.999$ and decreased by 0.001 until the
208 algorithm produced a PSD correlation matrix or $k = 0$.¹

209 **Knol-Berger Algorithm (KB; 1991).** In contrast to the Higham (2002) and
210 Bentler-Yuan (2011) smoothing algorithms, the Knol-Berger algorithm is a non-iterative
211 procedure in which the negative eigenvalues of \mathbf{R}_- are replaced with some small positive
212 value. The first step in the Knol-Berger algorithm is to compute the eigendecomposition of
213 the $p \times p$ indefinite correlation matrix, as defined in the previous section. Next, a matrix Λ_+
214 is created by setting all negative elements of Λ_- equal to some user-specified small, positive
215 constant. Finally, a smoothed, PSD correlation matrix, \mathbf{R}_{KB} , is constructed by replacing Λ_-
216 with Λ_+ in the eigendecomposition of \mathbf{R}_- and then scaling to ensure a unit diagonal and
217 that the absolute value of all off-diagonal elements is less than or equal to one:

$$\mathbf{R}_{\text{KB}} = [\text{dg}(\mathbf{V}\Lambda_+\mathbf{V}')]^{-1/2}\mathbf{V}\Lambda_+\mathbf{V}'[\text{dg}(\mathbf{V}\Lambda_+\mathbf{V}')]^{-1/2}, \quad (5)$$

218 where the dg operator is defined such that $\text{dg}(\mathbf{A})$ returns a diagonal matrix containing
219 the diagonal elements of \mathbf{A} (Magnus & Neudecker, 2019, p. 6).

¹ Bentler and Yuan suggest using $k = 0.96$ (Bentler & Yuan, 2011, p. 120) but suggest that the precise value of k does not matter a great deal as long as k is marginally less than one.

220 **The Common Factor Model**

221 The linear factor analysis model is used to describe the variance of each observed
222 variable in terms of the contributions of a small number of latent common factors and a
223 specific factor unique to that variable (Wirth & Edwards, 2007). In the common factor
224 model, the population correlation matrix, \mathbf{P} , can be expressed as:

$$\mathbf{P} = \mathbf{F}\Phi\mathbf{F}' + \boldsymbol{\Theta}^2, \quad (6)$$

225 where \mathbf{P} is a $p \times p$ population correlation matrix for p observed variables, \mathbf{F} is a $p \times m$
226 factor loading matrix for m common factors, Φ is an $m \times m$ matrix of correlations between
227 the m common factors, and $\boldsymbol{\Theta}^2$ is a $p \times p$ diagonal matrix containing the unique variances.

228 Although the common factor analysis model represented in Equation (6) is often useful,
229 many authors have remarked that it constitutes an oversimplification of the complex
230 processes that generate real, observed data (Cudeck & Henly, 1991; MacCallum & Tucker,
231 1991; MacCallum et al., 2001). Tucker et al. (1969) suggested that the lack-of-fit between
232 the common factor model and the complex processes underlying real data could be
233 represented by modeling a large number of minor common factors of small effect. The model
234 Tucker et al. (1969) proposed can be written as:

$$\mathbf{P} = \mathbf{F}\Phi\mathbf{F}' + \boldsymbol{\Theta}^2 + \mathbf{WW}', \quad (7)$$

235 where \mathbf{W} is a $p \times q$ matrix containing factor loadings for the $q \gg m$ minor factors (Briggs &
236 MacCallum, 2003, p. 32). Given our expectation that the common factor model is not a
237 perfect representation of any real-world data-generating process we might wish to represent,
238 Equation (7) is arguably preferable to Equation (6) for simulating realistic data (Briggs &
239 MacCallum, 2003; Hong, 1999).

²⁴⁰ **Factor Extraction Methods**

²⁴¹ Various factor extraction methods have been proposed for estimating item factor
²⁴² loadings, factor correlations, and unique item variances. One purpose of this study was to
²⁴³ determine whether the effects of matrix smoothing method on factor loading recovery differ
²⁴⁴ depending on which factor extraction method is used. To that end, three of the most
²⁴⁵ commonly-used factor extraction methods (Fabrigar, Wegener, MacCallum, & Strahan, 1999)
²⁴⁶ were used in the present simulation study: principal axis (PA), ordinary least-squares (OLS),
²⁴⁷ and maximum-likelihood (ML).

²⁴⁸ **Principal Axis Factor Analysis.** Principal axis (PA) factor analysis is
²⁴⁹ conceptually similar to principal components analysis (PCA). Whereas PCA seeks to find a
²⁵⁰ low-dimensional approximation of the full observed correlation matrix, PA seeks to find a
²⁵¹ low-dimensional approximation of the reduced correlation matrix, \mathbf{R}_* (i.e., the observed
²⁵² correlation matrix, \mathbf{R} , with communalities on the diagonal). Because the true communalities
²⁵³ are unknown, principal axis factor analysis starts by using estimated communalities to form
²⁵⁴ \mathbf{R}_* .² The eigenvalues of \mathbf{R}_* are then taken to be the updated communality estimates. These
²⁵⁵ updated estimates replace the previous estimates on the diagonal of \mathbf{R}_* and the procedure
²⁵⁶ iterates until the sum of the differences between the communality estimates from the current
²⁵⁷ and previous iterations is less than some small convergence criterion.

²⁵⁸ **Ordinary Least-Squares Factor Analysis.** The ordinary least-squares factor
²⁵⁹ analysis method (OLS; also known as “minres”; Comrey, 1962) seeks to minimize the sum of
²⁶⁰ squared differences between the sample correlation matrix, \mathbf{R} , and $\hat{\mathbf{P}} = \hat{\mathbf{F}}\hat{\Phi}\hat{\mathbf{F}}' + \hat{\Theta}^2$, the
²⁶¹ correlation matrix implied by the estimated factor model corresponding to Equation (6).

² Many methods of estimating communalities have been proposed, the most common of which are the squared multiple correlation between each variable and the other variables (Dwyer, 1939; Mulaik, 2009, p. 182; Roff, 1936) and the maximum absolute correlation between each variable and the other variables (Mulaik, 2009, p. 175; Thurstone, 1947). However, the particular choice of initial communality estimates has been shown to not have a large effect on the final solution when the convergence criterion is sufficiently stringent (Widaman & Herringer, 1985).

²⁶² The OLS discrepancy function can then be written as

$$F_{OLS}(\mathbf{R}, \hat{\mathbf{P}}) = \frac{1}{2} \text{tr} [(\mathbf{R} - \hat{\mathbf{P}})^2], \quad (8)$$

²⁶³ where tr is the trace operator (Magnus & Neudecker, 2019, p. 11) and $\text{tr} [(\mathbf{R} - \hat{\mathbf{P}})^2]$ is the
²⁶⁴ trace (sum of the diagonal elements) of the matrix formed by $(\mathbf{R} - \hat{\mathbf{P}})^2$. OLS does not give
²⁶⁵ additional weight to residuals corresponding to large correlations and requires no
²⁶⁶ assumptions about the population distributions of the variables (Briggs & MacCallum, 2003).

²⁶⁷ **Maximum-Likelihood Factor Analysis.** The maximum likelihood factor analysis
²⁶⁸ algorithm (ML) is similar to OLS in that it seeks to minimize the discrepancy between \mathbf{R}
²⁶⁹ and $\hat{\mathbf{P}}$. Unlike OLS, however, ML assumes that all variables are multivariate normal in the
²⁷⁰ population. Then, we can write the discrepancy function to be minimized as an alternative
²⁷¹ form of the multivariate normal log-likelihood function,

$$F_{ML}(\mathbf{R}, \hat{\mathbf{P}}) = \log |\hat{\mathbf{P}}| - \log |\mathbf{R}| + \text{tr}(\mathbf{S}\hat{\mathbf{P}}^{-1}) - p. \quad (9)$$

²⁷² In addition to the distributional assumptions required by ML factor analysis, the method
²⁷³ also assumes that the only source of error in the model is sampling error. Consequently,
²⁷⁴ large correlations (having relatively small standard errors) are fit more closely than small
²⁷⁵ correlations (with relatively large standard errors) under maximum likelihood factor analysis
²⁷⁶ (Briggs & MacCallum, 2003). Also note that when \mathbf{R} is indefinite, $|\mathbf{R}|$ is negative and
²⁷⁷ $\log |\mathbf{R}|$ is undefined. Therefore, indefinite covariance or correlation matrices cannot be used
²⁷⁸ as input for maximum likelihood factor analysis.

²⁷⁹ Simulation Procedure

²⁸⁰ I conducted a simulation study to evaluate four approaches to dealing with indefinite
²⁸¹ tetrachoric correlation matrices (applying matrix smoothing using the Higham [2002],
²⁸² Bentler-Yuan [2011], or Knol-Berger [1991] algorithms, or leaving indefinite tetrachoric

283 matrices unsmoothed) in the context of exploratory factor analysis. The simulation study
284 was designed to address two primary questions. First, which smoothing method (Higham,
285 Bentler-Yuan, Knol-Berger, or None) produced (possibly) smoothed correlation matrices
286 (\mathbf{R}_{Sm}) that most closely approximated the corresponding population correlation matrices
287 (\mathbf{R}_{Pop})? Second, which smoothing method produced correlation matrices that led to the best
288 estimates of the population factor loading matrix when used in exploratory factor analyses?

289 In the first step of the simulation study, I generated random sets of binary data from a
290 variety of orthogonal factor models with varying numbers of major common factors
291 (Factors $\in \{1, 3, 5, 10\}$). Using the method of Tucker et al. (1969), I also incorporated the
292 effects of model approximation error into the data by including 150 minor common factors in
293 each population model. In total, these 150 minor common factors accounted for 0%, 10%, or
294 30% ($v_E \in \{0, .1, .3\}$) of the uniqueness variance of the error-free model (i.e., the model with
295 only the major common factors). These conditions were chosen to represent models with
296 perfect, good, or moderate model fit, resembling the conditions used by Briggs and
297 MacCallum (2003). These three levels of model error variance ensured that both ideal
298 ($v_E = 0$) and more empirically-plausible levels of model error variance ($v_E \in \{.1, .3\}$) were
299 considered in this study.

300 In addition to systematically varying the number of major factors and the proportion
301 of uniqueness variance accounted for by model approximation error, I also varied the number
302 of factor indicators (i.e., items loading on each factor; Items/Factor $\in \{5, 10\}$), and the
303 number of subjects per item (Subjects/Item $\in \{5, 10, 15\}$). The total numbers of items (p)
304 and sample sizes (N) for each factor number condition can be found in Table 1. Each item
305 loaded on only one factor and item factor loadings were uniformly fixed at one of three levels
306 (Loading $\in \{.3, .5, .8\}$). Though “rules-of-thumb” for factor loadings vary, Hair, Black, Babin,
307 and Anderson (2018, p. 151) suggest that “[f]actor loadings in the range of ± 0.30 to ± 0.40
308 are considered to meet the minimal level for interpretation of structure”, and “[l]oadings

³⁰⁹ ± 0.50 or greater are considered practically significant.” Moreover, factor loadings of ± 0.8 are
³¹⁰ considered to be high (MacCallum et al., 2001). Thus, the three factor loadings investigated
³¹¹ in this study were chosen to represent low, moderate, and high levels of factor salience.

³¹² The combinations of the independent variables specified above resulted in a
³¹³ fully-crossed design with 4 (Factors) \times 3 (Model Error, v_E) \times 2 (Items/Factor) \times 3
³¹⁴ (Subjects/Item) \times 3 (Loading) = 216 unique conditions. For each of these conditions, the
³¹⁵ `simFA()` function in the R (Version 3.6.2; R Core Team, 2019)³ *fungible* package (Version
³¹⁶ 1.95.4.8; Waller, 2019) was used to generate 1,000 random sets of binary data.

³¹⁷ Data Generation

³¹⁸ Each data set in the simulation was generated as follows. First, a model-implied
³¹⁹ population correlation matrix, \mathbf{R}_{Pop} , was generated using

$$\mathbf{R}_{\text{Pop}} = \mathbf{F}\Phi\mathbf{F}' + \boldsymbol{\Theta}^2 + \mathbf{W}\mathbf{W}'. \quad (10)$$

³²⁰ Here, \mathbf{F} denotes a $p \times m$ matrix of major factor loadings with simple structure such that
³²¹ each factor had exactly p/m salient loadings (fixed at the value indicated by the level of
³²² Loading) and all other loadings fixed at zero. Because only orthogonal models were
³²³ considered in this study, the factor correlation matrix Φ was an $m \times m$ identity matrix.

³ Additionally, I used the following R packages: *arm* (Version 1.10.1; Gelman & Su, 2018), *broom.mixed* (Version 0.2.4; Bolker & Robinson, 2019), *car* (Version 3.0.7; Fox & Weisberg, 2019), *dplyr* (Version 0.8.5; Wickham et al., 2019), *forcats* (Version 0.5.0; Wickham, 2019a), *ggplot2* (Version 3.3.0; Wickham, 2016), *here* (Version 0.1.11; Müller, 2017), *knitr* (Version 1.28; Xie, 2015), *koRpus* (Version 0.11.5; Michalke, 2018a, 2019), *koRpus.lang.en* (Version 0.1.3; Michalke, 2019), *latex2exp* (Version 0.4.0; Meschiari, 2015), *lattice* (Version 0.20.38; Sarkar, 2008), *lme4* (Version 1.1.23; Bates, Mächler, Bolker, & Walker, 2015), *MASS* (Version 7.3.51.4; Venables & Ripley, 2002), *Matrix* (Version 1.2.18; Bates & Maechler, 2019), *merTools* (Version 0.5.0; Knowles & Frederick, 2019), *papaja* (Version 0.1.0.9942; Aust & Barth, 2018), *patchwork* (Version 1.0.0; Pedersen, 2019), *purrr* (Version 0.3.4; Henry & Wickham, 2019), *questionr* (Version 0.7.0; Barnier, Briatte, & Larmarange, 2018), *readr* (Version 1.3.1; Wickham, Hester, & Francois, 2018), *sfsmisc* (Version 1.1.4; Maechler, 2019), *stringr* (Version 1.4.0; Wickham, 2019b), *syly* (Version 0.1.5; Michalke, 2018b), *texreg* (Version 1.36.23; Leifeld, 2013), *tibble* (Version 3.0.1; Müller & Wickham, 2019), *tidyR* (Version 1.0.2.9000; Wickham & Henry, 2019), *tidyverse* (Version 1.3.0; Wickham, Averick, et al., 2019), *viridis* (Version 0.5.1; Garnier, 2018), and *wordcountaddin* (Version 0.3.0.9000; Marwick, 2019).

The $p \times q$ matrix of minor common factor loadings, \mathbf{W} , was constructed in multiple steps. First, a $p \times q$ provisional matrix, \mathbf{W}^* , was generated such that the i th column of \mathbf{W}^* consisted of p independent samples from $\mathcal{N}(0, (1 - \epsilon)^{2(i-1)})$ where $\epsilon \in [0, 1]$ was a user-specified constant. The value of ϵ determined how the minor common factor (error) variance was distributed. Values of ϵ close to zero resulted in the error variance being spread relatively equally among the minor common factors. Values of ϵ close to one resulted in error variance primarily being distributed to the first minor factor, with the remaining variance distributed to the other minor factors in a decreasing geometric sequence. To ensure that the minor common factors accounted for the specified proportion of uniqueness variance (denoted as v_E), \mathbf{W}^* was scaled to create \mathbf{W} . This scaling was done in several steps. First, a diagonal matrix $\Theta_{p \times p}^*$ was created such that

$$\Theta^* = \mathbf{I}_p - \text{dg}(\mathbf{F}\mathbf{F}'), \quad (11)$$

where $\text{dg}(\mathbf{F}\mathbf{F}')$ is to be read as the diagonal matrix formed from the diagonal entries in $\mathbf{F}\mathbf{F}'$ and \mathbf{I}_p denotes a $p \times p$ identity matrix. Then the matrix \mathbf{W} was formed using

$$\mathbf{W} = (\text{dg}(\mathbf{W}^*\mathbf{W}^{*\prime})^{-1}\Theta^*v_E)^{1/2}\mathbf{W}^*. \quad (12)$$

This process ensured that the q minor common factors accounted for the specified proportion of the variance not accounted for by the major common factors. The \mathbf{W} matrix was then used to create the diagonal matrix of unique variances, $\Theta^2 = \mathbf{I}_p - \text{Diag}(\mathbf{F}\mathbf{F}' + \mathbf{W}\mathbf{W}')$. The \mathbf{F} , Θ^2 , and \mathbf{W} matrices were then used to construct population correlation matrix, \mathbf{R}_{Pop} , as shown in Equation (10).

Having specified the elements of the population common factor model, the next step in the data-generation procedure was to draw a sample correlation matrix, \mathbf{R} , (for a given sample size, N) from \mathbf{R}_{Pop} using the method of Kshirsagar (1959; see also, Browne, 1968).

345 The sample correlation matrix was then used to generate a matrix of continuous data,
 346 $\mathbf{X}_{N \times p} = (X_1, \dots, X_N)'$, where $X \sim \mathcal{N}_p(\mathbf{0}_p, \mathbf{R})$. To obtain binary responses from the
 347 continuous data, items were assigned classical item difficulties (d ; i.e., the expected
 348 proportion of correct responses, Crocker & Algina, 1986) at equal intervals between 0.15 and
 349 0.85. For example, items in a five-item data set were assigned classical item difficulties of
 350 .150, .325, .500, .675, and .850. The classical item difficulties were used to obtain threshold
 351 values, t , such that $1 - \Phi(t) = d$. Using these thresholds, the continuous data were converted
 352 to binary data. If a data set had any homogeneous item response vectors (i.e., had one or
 353 more items with zero variance), the data set was discarded and a new sample of data was
 354 generated until all items had non-homogeneous response vectors. Homogeneous response
 355 vectors were not allowed because such response vectors can lead to poorly-estimated
 356 tetrachoric correlations (Brown & Benedetti, 1977).

357 Next, a tetrachoric correlation matrix was computed for each simulated binary data set.
 358 Tetrachoric correlation matrices were calculated using the `tetcor()` function in the R
 359 *fungible* package (Waller, 2019), which computes maximum likelihood tetrachoric correlation
 360 coefficients (Brown & Benedetti, 1977; Olsson, 1979). If a tetrachoric correlation matrix was
 361 indefinite, the Higham (2002), Bentler-Yuan (2011), and Knol-Berger (1991) matrix
 362 smoothing algorithms were applied to the indefinite tetrachoric correlation matrix to produce
 363 three smoothed, PSD correlation matrices. Matrix smoothing was done using the
 364 `smoothAPA()`, `smoothBY()`, and `smoothKB()` implementations of the Higham (2002),
 365 Bentler-Yuan (2011), and Knol-Berger (1991) algorithms in the *fungible* package.

366 In the final step of the simulation procedure, three exploratory factor analysis
 367 algorithms (principal axis [PA], ordinary least squares [OLS], and maximum likelihood [ML])
 368 were applied to each of the indefinite tetrachoric correlation matrices and the PSD,
 369 smoothed correlation matrices. Because ML does not work with indefinite correlation or
 370 covariance matrices as input, ML was conducted on the Pearson correlation matrix (rather

371 than the indefinite tetrachoric correlation matrix) when no smoothing was applied. Each of
 372 the factor solutions were rotated using a quartimin rotation (Carroll, 1957; Jennrich, 2002)
 373 and aligned to match the corresponding population factor loading matrix such that the least
 374 squares discrepancy between the matrices was minimized. The alignment step ensured that
 375 the elements of each estimated factor loading matrix were matched (in order and sign) to the
 376 elements of the corresponding population factor loading matrix. These rotation and
 377 alignment steps were accomplished using the `faMain()` and `faAlign()` functions in the R
 378 *fungible* package (Waller, 2019). Code for all aspects of this study is available at
 379 https://github.umn.edu/krach018/masters_thesis.

380

Results

381 Recovery of \mathbf{R}_{Pop}

382 One of the primary reasons for conducting the present simulation study was to
 383 determine which of the three investigated smoothing methods—the Higham (2002),
 384 Bentler-Yuan (2011), or Knol-Berger (1991) algorithms—resulted in smoothed correlation
 385 matrices that were closest to the correlation matrix implied by the major factor model (i.e.,
 386 the factor model not including the minor factors). In particular, I examined whether
 387 smoothed correlation matrices were closer to the model-implied correlation matrix than the
 388 unsmoothed, indefinite correlation matrix. In this context, the scaled distance between two
 389 $p \times p$ correlation matrices $\mathbf{A} = \{a_{ij}\}$ and $\mathbf{B} = \{b_{ij}\}$ was computed as:

$$D_s(\mathbf{A}, \mathbf{B}) = \sqrt{\sum_{i=1}^{p-1} \sum_{j=i+1}^p \frac{(a_{ij} - b_{ij})^2}{p(p-1)/2}}. \quad (13)$$

390 To understand which of the smoothing algorithms most often produced a smoothed
 391 correlation matrix, \mathbf{R}_{Sm} , that was closest to the model-implied correlation matrix, \mathbf{R}_{Pop} , I
 392 calculated $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$ for each \mathbf{R}_{Sm} obtained from the 124,346 indefinite tetrachoric

correlation matrices.⁴ Small values of $D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})$ indicated that the smoothed correlation matrix was a good approximation of \mathbf{R}_{Pop} , whereas large values indicated that \mathbf{R}_{Sm} was a poor approximation of \mathbf{R}_{Pop} . After excluding three cases where the Higham (2002) algorithm failed to converge, I fit a linear mixed-effects model (Model 1A) regressing $\log D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})$ on all of the simulation design variables and their two-way interactions. Additionally, a random intercept was estimated for every unique indefinite correlation matrix to account for the correlation between observations corresponding the same indefinite correlation matrix.⁵

The estimated fixed-effect coefficients are shown in Figure 1. A full summary table for the model is contained in Table 2. Figure 1 and Table 2 also summarize the results of a second model (Model 1B) that included second-degree polynomial terms for number of factors, factor loading, and subjects per item in addition to the terms included in Model 1A. The results in Table 2 indicated that Model 1B should be preferred based on the AIC (Akaike, 1973) and BIC (e.g., Hastie, Tibshirani, & Friedman, 2009) criteria. Therefore, coefficient estimates and estimated marginal means reported in this section were obtained using Model 1B.

The design variables that most influenced $D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})$ values can be seen in Figure 1, which shows coefficient estimates with 99% confidence intervals, ordered by size. Note that exponentiated coefficients less than 1.01 and greater than 0.99 were omitted from the figure to conserve space. Figure 1 shows that only a few variables had non-trivial effects on population matrix recovery. In particular, the three largest effects were for number of factors ($\hat{b} = -0.52$, $SE = 0.00$, $e^{-0.52} = 0.59$), number of items per factor ($\hat{b} = -0.26$, $SE = 0.00$, $e^{-0.26} = 0.77$), and number of subjects per item ($\hat{b} = -0.25$, $SE = 0.00$, $e^{-0.25} = 0.78$). These estimated effects were all negative, indicating better recovery of the population correlation matrix for models with larger numbers of major factors, larger numbers of items per factor, and larger numbers of subjects per item. The effects of number of factors and

⁴ A table reporting the percent of tetrachoric correlation matrices that were indefinite can be found in Appendix B.

⁵ All numeric predictors were scaled to have a mean of zero and variance of one prior to analysis. Diagnostic plots are shown in Appendix A.

number of subjects per item were somewhat offset, however, by large (positive) estimated effects for the squared number of factors ($\hat{b} = 0.19$, $SE = 0.00$, $e^{0.19} = 1.21$) and squared number of subjects per item ($\hat{b} = 0.08$, $SE = 0.00$, $e^{0.08} = 1.08$) terms. The effects of all of the independent variables can be more easily understood by looking at Figure 2, which shows estimated marginal mean $D_s(\mathbf{R}_{\text{Pop}}, \mathbf{R}_{\text{Sm}})$ values (and 99% confidence intervals) at each level of number of factors, number of subjects per item, number of items per factor, factor loading, and smoothing method.⁶

The effects most relevant to the research question were the effects of the smoothing methods and their interactions with other variables. These effects were all relatively small, but can still be seen in Figure 2. For instance, the Bentler-Yuan algorithm (2011) had the largest (negative) main effect ($\hat{b} = -0.06$, $SE = 0.00$, $e^{-0.06} = 0.94$), closely followed by the Knol-Berger (1991; $\hat{b} = -0.01$, $SE = 0.00$, $e^{-0.01} = 0.99$) and Higham (2002; $\hat{b} = -0.01$, $SE = 0.00$, $e^{-0.01} = 0.99$) algorithms. These results suggest that all three algorithms generally led to smoothed correlation matrices that were closer to their population counterparts than were the unsmoothed, indefinite correlation matrices. However, the differences among the smoothing algorithms were largest for conditions with small numbers of subjects per item, small numbers of items per factor, and low factor loadings, as shown in Figure 2. Indeed, the results show that the application of matrix smoothing was most beneficial in conditions where \mathbf{R}_{Pop} was poorly estimated, regardless of which smoothing algorithm was used (or whether matrix smoothing was applied at all). In conditions where \mathbf{R}_{Pop} tended to be recovered better overall, there were at best only small differences between the four smoothing methods.

439

⁶ Estimated marginal means were used to summarize results because the data were unbalanced (due to only using indefinite tetrachoric correlation matrices in the analyses). Additional tables and figures showing results from the raw data can be found in Appendix B.

440 **Recovery of Factor Loadings**

441 I next analyzed the simulation results in terms of factor loading recovery. In particular,
 442 I was interested in whether factor analysis of smoothed indefinite correlation matrices led to
 443 better factor loading estimates compared to when factor analysis was conducted on the
 444 indefinite correlation matrices directly. I was also interested in whether particular smoothing
 445 methods led to better factor loading estimates than others and whether the interactions
 446 between smoothing methods and the other variables (e.g., number of items per factor,
 447 number of subjects per item, factor analysis method, etc.) affected factor loading estimation.
 448 For the purposes of these analyses, I evaluated factor loading recovery using the
 449 root-mean-square error (RMSE) between the estimated and population factor loadings for
 450 the major factors. Given a matrix of estimated major factor loadings $\hat{\mathbf{F}} = \{\hat{f}_{ij}\}_{p \times m}$, and the
 451 corresponding matrix of population major factor loadings, $\mathbf{F} = \{f_{ij}\}_{p \times m}$,

$$\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}}) = \sqrt{\sum_{i=1}^p \sum_{j=1}^m \frac{(f_{ij} - \hat{f}_{ij})^2}{pm}}. \quad (14)$$

452 To determine which smoothing method resulted in the best factor loading estimates, I
 453 calculated the $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ for each pair of estimated and population factor loading
 454 matrices corresponding to the (possibly) smoothed indefinite tetrachoric correlation matrices.
 455 Relatively small $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values indicated that the estimated factor loading matrices
 456 were more similar to their corresponding population factor loading matrices, whereas larger
 457 $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values indicated poorly-estimated factor loading matrices. As in the previous
 458 section, the four cases where the Higham (2002) algorithm did not converge were not
 459 included in my analyses. Furthermore, cases where PA failed to converge were also not
 460 included. In total, there were 2,714 cases where the PA algorithm did not converge
 461 (convergence rate = 99.5%) and only four cases where the ML algorithm did not converge
 462 (convergence rate > 99.9%).

463 Using the converged data, I fit a mixed-effects model (Model 2A) regressing
464 $\log \text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ on number of subjects per item, number of items per factor, number of
465 factors, factor loading, model error, smoothing algorithm, factor analysis method (PA, OLS,
466 or ML), all two-way interactions between these variables, and a random intercept estimated
467 for every unique indefinite correlation matrix.⁷ I also fit a second mixed-effects model
468 (Model 2B) with additional second-degree polynomial terms for number of subjects per item,
469 number of factors, and factor loading. The results for both models are summarized in Table
470 3 and indicated that Model 2B should be preferred to Model 2A based on the AIC and BIC
471 criteria. Therefore, coefficient estimates and estimated marginal means reported in this
472 section were obtained using Model 2B.

473 To show the variables that most affected $\log \text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$, ordered, exponentiated
474 coefficient estimates with 99% confidence intervals for Model 2B are shown in Figure 3 (note
475 that exponentiated coefficients less than 1.01 and greater than 0.99 were omitted to conserve
476 space). Figure 3 shows that factor loading, items per factor, number of factors, model error,
477 and subjects per item had relatively large effects on $\log \text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$. Additionally, many of
478 the polynomial terms and interactions between these variables also had relatively large
479 estimated effects (see Figure 3 and Table 3). To better understand the effects of each of
480 these design variables on $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$, Figure 4 shows estimated marginal means
481 conditioned on number of factors, model error, number of subjects per item, number of items
482 per factor, smoothing method, factor extraction method, and factor loading. As in the
483 previous section, estimated marginal means were used instead of raw means because each
484 condition of the design had a different number of observations due to only using results for
485 indefinite tetrachoric correlation matrices.

486 Concerning the primary question of interest in this section, the coefficient estimates

⁷ All numeric predictors were scaled to have a mean of zero and variance of one prior to analysis. Diagnostic plots can be found in Appendix A.

487 from Model 2B (and the marginal means shown in Figure 4) indicated that choice of
 488 smoothing method generally had little impact on RMSE($\mathbf{F}, \hat{\mathbf{F}}$). Of the effects involving
 489 smoothing methods, only the effects involving the Bentler-Yuan algorithm were
 490 non-negligible. Therefore, only the application of the Bentler-Yuan algorithm to indefinite
 491 tetrachoric correlation matrices seemed to be related to any improvement (or indeed,
 492 difference) in RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values when used for factor analysis. Moreover, even the
 493 estimated main effect associated with the Bentler-Yuan algorithm (2011) was quite modest
 494 ($\hat{b} = -0.03, SE = 0.00, e^{-0.03} = 0.97$) and was offset by the positive estimated interaction
 495 effects with factor loading ($\hat{b} = 0.02, SE = 0.00, e^{0.02} = 1.02$) and items per factor ($\hat{b} = 0.02,$
 496 $SE = 0.00, e^{0.02} = 1.02$). These effects are evident in Figure 4, which shows that the
 497 Bentler-Yuan algorithm (2011) led to slightly lower RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values for conditions with
 498 low factor loadings and few items per factor, but led to nearly identical results to the
 499 alternative methods as factor loading magnitude and number of items per factor increased.

500 Although choice of smoothing method did not have a large influence on factor loading
 501 recovery, the other design variables did have an influence on RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values. The effects
 502 of these other variables (including interactions and polynomial terms) are best understood
 503 using the marginal means shown in Figure 4. Considering first the effect of factor loading,
 504 Figure 4 shows that RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values tended to decrease as factor loadings increased
 505 ($\hat{b} = -0.57, SE = 0.00, e^{-0.57} = 0.57$). Interestingly, there was also a relatively large
 506 interaction between factor loading and ML factor extraction ($\hat{b} = 0.22, SE = 0.00,$
 507 $e^{0.22} = 1.25$) such that ML seemed to benefit less than PA or OLS from higher factor
 508 loadings.

509 After factor loading, the number of items per factor had the largest estimated effect on
 510 RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values ($\hat{b} = -0.37, SE = 0.00, e^{-0.37} = 0.69$). As can be seen in Figure 4,
 511 increasing the number of items per factor tended to improve factor loading recovery (i.e., led
 512 to lower RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values). There was a similar (albeit smaller) effect for the number of

513 subjects per item, such that increasing the number of subjects per item tended to lead to
 514 smaller RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values ($\hat{b} = -0.18, SE = 0.00, e^{-0.18} = 0.84$). However, both the effects
 515 of number of items per factor and number of subjects per item were affected by their
 516 associated polynomial terms and interactions. For instance, the effect of increasing the
 517 number of items per factor was largest when factor loadings were relatively small, as
 518 indicated by the interaction between the number of items per factor and squared factor
 519 loadings ($\hat{b} = 0.14, SE = 0.00, e^{0.14} = 1.15$). The effect of the number of items per factor
 520 was also influenced by model error (as will be discussed shortly). Concerning the effect of the
 521 number of subjects per item, there was a large estimated effect for squared number of
 522 subjects per item such that the beneficial effect of increasing the number of subjects
 523 dissipated as the number of subjects per item increased ($\hat{b} = 0.16, SE = 0.00, e^{0.16} = 1.17$).

524 Concerning the effect of number of factors on RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values, Figure 4 shows that
 525 RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values tended to be lower for conditions with more factors compared to those
 526 with fewer factors ($\hat{b} = -0.34, SE = 0.00, e^{-0.34} = 0.71$). Moreover, the decrease in
 527 RMSE($\mathbf{F}, \hat{\mathbf{F}}$) seemed to be nonlinear, such that the effect of increasing the number of factors
 528 was largest when there were relatively few factors, as indicated by the quadratic term for
 529 number of factors ($\hat{b} = 0.09, SE = 0.00, e^{0.09} = 1.09$). On its face, these effects seem to
 530 suggest that models with large numbers of major factors led to better factor loading recovery
 531 than those with fewer factors. However, another explanation for this effect involves the total
 532 numbers of subjects and items. Whereas number of items per factor and number of subjects
 533 per item were fully-crossed with number of factors, the total sample size and total number of
 534 items for each data set were confounded with number of factors. In other words, conditions
 535 with larger numbers of factors tended to include more total subjects and items. The strong
 536 relationship between log RMSE($\mathbf{F}, \hat{\mathbf{F}}$) and sample size can be clearly seen in Figure 5, which
 537 shows that log RMSE($\mathbf{F}, \hat{\mathbf{F}}$) decreased as sample size increased. Therefore, it seems
 538 reasonable that the effect of number of factors might be better understood as being related
 539 to the total number of items and subjects in a data set. Similarly, the negative interaction

540 between number of factors and ML ($\hat{b} = -0.20$, $SE = 0.00$, $e^{-0.20} = 0.82$) could be
 541 interpreted instead as an interaction between total number of items or subjects and ML.

542 Moving next to model error, Model 2B indicated that increasing the proportion of
 543 uniqueness variance accounted for by the minor factors (v_E) was associated with worse factor
 544 loading recovery ($\hat{b} = 0.16$, $SE = 0.00$, $e^{0.16} = 1.17$). Additionally, the detrimental effect of
 545 model error on factor loading recovery seemed to worsen as the number of factors and
 546 number of items per factor increased (see Figure 4 and Table 3). On the other hand, model
 547 error seemed to have less of an impact on RMSE($\mathbf{F}, \hat{\mathbf{F}}$) as factor loadings increased, as can
 548 also be seen in Figure 4. A potential explanation for this effect is that model error accounted
 549 for less of the total variance in conditions with high factor loadings because the levels of
 550 model error were defined as proportions of the uniqueness variance. I.e., conditions with high
 551 factor loadings had small uniqueness variances and correspondingly small model error
 552 variances.

553 Another notable effect involving model approximation error was the interaction
 554 between model error and ML factor extraction ($\hat{b} = -0.03$, $SE = 0.00$, $e^{-0.03} = 0.97$). This
 555 result indicated that, of the three factor extraction methods, ML was less affected by model
 556 error than OLS or PA. Moreover, the main effect of ML indicated that it led to lower overall
 557 RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values than OLS or PA when all other variables were held constant ($\hat{b} = -0.05$,
 558 $SE = 0.00$, $e^{-0.05} = 0.95$). However, the previously-discussed interactions between ML and
 559 factor loading, number of items per factor, and number of subjects per item indicated that
 560 ML led to better results than PA or OLS only when the numbers of subjects per item and
 561 items per factor were small and factor loadings were low. In conditions with higher numbers
 562 of subjects and items, OLS and PA led to better (and highly similar) results in terms of
 563 RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values.

565

Discussion

566 The current study examined how the application of three matrix smoothing algorithms
567 (the Higham [2002], Bentler-Yuan [2011], and Knol-Berger [1991] algorithms) to indefinite
568 tetrachoric correlation matrices affected both (a) the recovery of the model-implied
569 population correlation matrix (\mathbf{R}_{Pop}), and (b) the recovery of the population item factor
570 loadings in EFA (compared to leaving the indefinite correlation matrices unsmoothed). With
571 respect to recovery of \mathbf{R}_{Pop} , I found that that three variables were most related to
572 $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$: (a) the number of major factors in the data-generating model, (b) the
573 number of subjects per item, and (c) the number of items per (major) factor. Increases in
574 any of these variables were associated with improved population correlation matrix recovery.
575 I also found that choice of smoothing method was somewhat related to population
576 correlation matrix recovery. Specifically, the application of any of the three investigated
577 matrix smoothing algorithms led to smoothed matrices were slightly closer to the population
578 correlation matrix (\mathbf{R}_{Pop}) than the unsmoothed, indefinite tetrachoric correlation matrices.
579 The results indicated that although the three smoothing algorithms led to very similar
580 $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$ values in most conditions, the Bentler-Yuan algorithm (2011) led to slightly
581 lower $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$ values in conditions with few subjects per item, few items per factor,
582 and low factor loadings.

583

Concerning factor loading recovery, the simulation study results indicated that choice
of smoothing algorithm—or, in fact, whether smoothing was applied at all—was not an
important determinant of factor loading recovery when EFA was applied to smoothed or
unsmoothed indefinite tetrachoric correlation matrices. Similar to the previous analyses, the
Bentler-Yuan algorithm (2011) led to slightly better results (i.e., lower $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values)
than the alternative smoothing methods when factor loadings were low and there were few
items per factor. Moreover, the Bentler-Yuan algorithm led to slightly better results when
paired with maximum likelihood factor extraction (ML) compared to when the ordinary least

squares (OLS) or principal axis (PA) extraction methods were used. However, the differences between the four smoothing methods (in terms of $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values) were never large enough to be of practical importance. Although smoothing method choice was not found to be important for determining factor loading recovery, many of the other design variables were found to be important. In particular, $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values were smallest for conditions with high factor loadings, many items per factor, and with little or no model approximation error. Moreover, the results indicated that the OLS and PA factor extraction methods led to highly similar results under all conditions. ML factor extraction method led to better results than OLS and PA in conditions with low factor loadings, few items per factor, and few subjects per item. The results also indicated that factor loading recovery for ML was less affected by model approximation error than were OLS or PA.

The results of this simulation study concerning both population correlation matrix recovery ($D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})$) and population factor loading recovery ($\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$) can be put in the context of previous research. First, the current results provided additional evidence that the application of matrix smoothing algorithms to indefinite tetrachoric correlation matrices led to, at most, only a small effect on factor loading estimates in subsequent factor analyses. This result lends additional support to the conclusion of Knol and Berger (1991) that the effect of applying matrix smoothing to indefinite tetrachoric correlation matrices prior to conducting factor analysis was negligible.

To the extent that there were small differences among the smoothing methods in terms of $D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})$ and $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$, the Bentler-Yuan algorithm (2011) tended to lead to slightly better results than the alternative algorithms. Although I am not aware of any previous comparisons of relative smoothing algorithm performance in terms of population correlation matrix or factor loading recovery, Debelak and Tran (2013) and Debelak and Tran (2016) both found that the Bentler-Yuan algorithm led to somewhat better results than the Higham (2002) or Knol-Berger (1991) algorithms used with indefinite polychoric

correlation matrices in the context of parallel analysis. These results, combined with the results from the present study, suggest that the Bentler-Yuan algorithm (2011) should be the default choice for smoothing indefinite tetrachoric or polychoric correlation matrices prior to conducting parallel analysis or factor analysis.

Limitations and Future Directions

As with any simulation study, the present simulation design was not able to cover the full range of realistic data scenarios. For instance, the simulation design included only orthogonal population factor models and did not allow for correlated factors. Moreover, the present study only included factor models with equal numbers of salient items per factor and fixed, uniform factor loadings. It might be the case that these loading matrices were overly-simplified and not representative of real data. Future research on this topic should investigate whether more complex factor loading and correlation structures affect the performance of matrix smoothing algorithms in terms of population correlation matrix recovery and factor loading recovery. Additionally, the present study only investigated the effects of matrix smoothing on indefinite tetrachoric correlation matrices. Further research should be done to investigate the effects of matrix smoothing on indefinite polychoric correlation matrices, as well as correlation matrices that are indefinite due to other causes such as indefinite correlation matrices calculated using pairwise deletion (Wothke, 1993) or composite correlation matrices used in meta-analysis (Furrow & Beretvas, 2005). Little is known about whether the mechanism or “cause” of indefinite correlation matrices affects their structure, or how these potential differences might interact with the application of matrix smoothing algorithms.

Future research should also investigate ways to side-step the problem of indefinite tetrachoric correlation matrices. For instance, Choi, Kim, Chen, and Dannels (2011) found that polychoric correlation matrices estimated using expected a posteriori (EAP) rather than maximum-likelihood estimation led to estimates that were negatively biased but produced

643 comparable (or smaller) RMSE values in terms of recovering the “true” correlations. It
644 seems plausible that the slight shrinkage induced by using EAP as an estimation method
645 would make indefinite tetrachoric or polychoric correlation matrices less common. Finally,
646 full-information maximum likelihood (FIML; Bock & Aitkin, 1981) can be used to estimate
647 model parameters directly and doesn’t require the estimation of a tetrachoric correlation
648 matrix. Future research should investigate whether the use of FIML (which is
649 computationally intensive, particularly with large models) offers any benefit in terms of
650 parameter recovery when applied to data sets corresponding to indefinite tetrachoric
651 correlation matrices.

652 Conclusion

653 Despite the lackluster improvement in factor loading recovery when factor analysis was
654 conducted on smoothed rather than indefinite tetrachoric correlation matrices, the
655 application of one of the three investigated matrix smoothing algorithms on indefinite
656 tetrachoric correlation matrices is still recommended. None of the smoothing algorithms
657 regularly led to worse results (in terms of factor loading recovery) compared to the
658 conditions where the indefinite correlation matrix was left unsmoothed. Moreover, all of the
659 smoothing algorithms investigated in this study are computationally inexpensive and are
660 readily available as functions in R packages. For instance, the *fungible* (Waller, 2019),
661 *sfsmisc* (Maechler, 2019), and *Matrix* (Bates & Maechler, 2019) packages all contain
662 implementations of at least one of the three smoothing algorithms discussed in this article.
663 In particular, the Bentler-Yuan algorithm (2011) often led to results that were at least as
664 good (and sometimes slightly better) than the alternative smoothing algorithms and
665 therefore seems a default choice of smoothing algorithm. Where the Bentler-Yuan algorithm
666 is not available, the Knol-Berger algorithm (1991) is an alternative that is fast, easily
667 implemented in most programming languages, does not have convergence issues, and
668 generally led to results comparable to the Bentler-Yuan algorithm.

These recommendations come with a strong caveat. Namely, no matrix smoothing algorithm can reasonably be considered a remedy or solution for indefinite tetrachoric correlation matrices. Instead, researchers should consider indefinite tetrachoric correlation matrices to be symptoms of larger problems (e.g., small sample sizes, bad items, etc.) and be aware that practical solutions such as gathering more data or discarding bad items are likely to lead to better results than the application of matrix smoothing algorithms. In particular, indefinite tetrachoric correlation matrices are less likely to occur when sample sizes are large relative to the number of items (see Table 1 in Debelak & Tran, 2013, p. 70), allowing researchers to avoid the question of how to properly deal with an indefinite tetrachoric correlation matrix entirely. If collecting more data is not possible, researchers should consider removing problematic items. In short, all three investigated smoothing algorithms are reasonable choices for dealing with indefinite tetrachoric correlation matrices prior to factor analysis and seem to offer a modest benefit (in terms of factor loading recovery) compared to leaving the indefinite tetrachoric correlation matrix unsmoothed. However, the application of these algorithms should be considered to be little more than a band-aid fix that does not address the underlying issues leading to indefinite tetrachoric correlation matrices nor to a marked improvement in factor loading recovery.

686

References

- 687 Akaike, H. (1973). *Information theory and an extension of the maximum likelihood principle*.
688 (B. N. Petrov & F. Caski, Eds.) (pp. 267–281). Budapest: Akademiai Kiado.
- 689 Aust, F., & Barth, M. (2018). *papaja: Create APA manuscripts with R Markdown*.
690 Retrieved from <https://github.com/crsh/papaja>
- 691 Banerjee, S., & Roy, A. (2014). *Linear algebra and matrix analysis for statistics*. Chapman;
692 Hall/CRC. <https://doi.org/10.1201/b17040>
- 693 Barnier, J., Briatte, F., & Larmarange, J. (2018). *Questionr: Functions to make surveys*
694 *processing easier*. Retrieved from <https://CRAN.R-project.org/package=questionr>
- 695 Bates, D., & Maechler, M. (2019). *Matrix: Sparse and dense matrix classes and methods*.
696 Retrieved from <https://CRAN.R-project.org/package=Matrix>
- 697 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models
698 using lme4. *Journal of Statistical Software*, 67(1), 1–48.
699 <https://doi.org/10.18637/jss.v067.i01>
- 700 Bentler, P. (1972). A lower-bound method for the dimension-free measurement of internal
701 consistency. *Social Science Research*, 1(4), 343–357.
702 [https://doi.org/10.1016/0049-089X\(72\)90082-8](https://doi.org/10.1016/0049-089X(72)90082-8)
- 703 Bentler, P., & Yuan, K.-H. (2011). Positive definiteness via off-diagonal scaling of a
704 symmetric indefinite matrix. *Psychometrika*, 76(1), 119–123.
705 <https://doi.org/10.1007/s11336-010-9191-3>
- 706 Birnbaum, A. L. (1968). Some latent trait models and their use in inferring an examinee's
707 ability. *Statistical Theories of Mental Test Scores*.

- 708 Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item
709 parameters: Application of an em algorithm. *Psychometrika*, 46(4), 443–459.
710 <https://doi.org/10.1007/BF02293801>
- 711 Bock, R. D., Gibbons, R., & Muraki, E. (1988). Full-information item factor analysis.
712 *Applied Psychological Measurement*, 12(3), 261–280.
713 <https://doi.org/10.1177/014662168801200305>
- 714 Bolker, B., & Robinson, D. (2019). *Broom.mixed: Tidying methods for mixed models*.
715 Retrieved from <https://CRAN.R-project.org/package=broom.mixed>
- 716 Box, G. E., & Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal
717 Statistical Society: Series B (Methodological)*, 26(2), 211–243.
718 <https://doi.org/10.1111/j.2517-6161.1964.tb00553.x>
- 719 Briggs, N. E., & MacCallum, R. C. (2003). Recovery of weak common factors by maximum
720 likelihood and ordinary least squares estimation. *Multivariate Behavioral Research*,
721 38(1), 25–56. https://doi.org/10.1207/S15327906MBR3801_2
- 722 Brown, M. B., & Benedetti, J. K. (1977). On the mean and variance of the tetrachoric
723 correlation coefficient. *Psychometrika*, 42(3), 347–355.
724 <https://doi.org/10.1007/BF02293655>
- 725 Browne, M. W. (1968). A comparison of factor analytic techniques. *Psychometrika*, 33(3),
726 267–334. <https://doi.org/10.1007/BF02289327>
- 727 Butler, S. M., & Louis, T. A. (1992). Random effects models with non-parametric priors.
728 *Statistics in Medicine*, 11(14-15), 1981–2000. <https://doi.org/10.1002/sim.4780111416>
- 729 Carroll, J. B. (1957). Biquartimin criterion for rotation to oblique simple structure in factor
730 analysis. *Science*, 126, 1114–1115. <https://doi.org/10.1126/science.126.3283.1114>

- 731 Choi, J., Kim, S., Chen, J., & Dannels, S. (2011). A comparison of maximum likelihood and
732 Bayesian estimation for polychoric correlation using monte carlo simulation. *Journal
733 of Educational and Behavioral Statistics*, 36(4), 523–549.
734 <https://doi.org/10.3102/1076998610381398>
- 735 Comrey, A. L. (1962). The minimum residual method of factor analysis. *Psychological
736 Reports*, 11(1), 15–18. <https://doi.org/10.2466/pr0.1962.11.1.15>
- 737 Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. Orlando,
738 FL: Wadsworth Publishing.
- 739 Cudeck, R., & Henly, S. J. (1991). Model selection in covariance structures analysis and the
740 “problem” of sample size: A clarification. *Psychological Bulletin*, 109(3), 512.
741 <https://doi.org/10.1037/0033-2909.109.3.512>
- 742 de Ayala, R. J. (2013). *The theory and practice of item response theory*. Guilford
743 Publications.
- 744 Debelak, R., & Tran, U. S. (2013). Principal component analysis of smoothed tetrachoric
745 correlation matrices as a measure of dimensionality. *Educational and Psychological
746 Measurement*, 73(1), 63–77. <https://doi.org/10.1177/0013164412457366>
- 747 Debelak, R., & Tran, U. S. (2016). Comparing the effects of different smoothing algorithms
748 on the assessment of dimensionality of ordered categorical items with parallel analysis.
749 *PLOS ONE*, 11(2), 1–18. <https://doi.org/10.1371/journal.pone.0148143>
- 750 Devlin, S. J., Gnanadesikan, R., & Kettenring, J. R. (1975). Robust estimation and outlier
751 detection with correlation coefficients. *Biometrika*, 62(3), 531–545.
752 <https://doi.org/10.1093/biomet/62.3.531>
- 753 Dillon, W. R., Kumar, A., & Mulani, N. (1987). Offending estimates in covariance structure

- 754 analysis: Comments on the causes of and solutions to Heywood cases. *Psychological*
755 *Bulletin*, 101(1), 126. <https://doi.org/10.1037/0033-2909.101.1.126>
- 756 Divgi, D. R. (1979). Calculation of the tetrachoric correlation coefficient. *Psychometrika*,
757 44(2), 169–172. <https://doi.org/10.1007/BF02293968>
- 758 Dong, H.-K. (1985). Non-Gramian and singular matrices in maximum likelihood factor
759 analysis. *Applied Psychological Measurement*, 9(4), 363–366.
760 <https://doi.org/10.1177/014662168500900404>
- 761 Dwyer, P. S. (1939). The contribution of an orthogonal multiple factor solution to multiple
762 correlation. *Psychometrika*, 4(2), 163–171. <https://doi.org/10.1007/BF02288494>
- 763 Fabrigar, L. R., Wegener, D. T., MacCallum, R. C., & Strahan, E. J. (1999). Evaluating the
764 use of exploratory factor analysis in psychological research. *Psychological Methods*,
765 4(3), 272. <https://doi.org/10.1037/1082-989X.4.3.272>
- 766 Fox, J., & Weisberg, S. (2019). *An R companion to applied regression* (Third). Thousand
767 Oaks CA: Sage. Retrieved from
768 <https://socialsciences.mcmaster.ca/jfox/Books/Companion/>
- 769 Furlow, C. F., & Beretvas, S. N. (2005). Meta-analytic methods of pooling correlation
770 matrices for structural equation modeling under different patterns of missing data.
771 *Psychological Methods*, 10(2), 227. <https://doi.org/10.1037/1082-989X.10.2.227>
- 772 Fushiki, T. (2009). Estimation of positive semidefinite correlation matrices by using convex
773 quadratic semidefinite programming. *Neural Computation*, 21(7), 2028–2048.
774 <https://doi.org/10.1162/neco.2009.04-08-765>
- 775 Garnier, S. (2018). *Viridis: Default color maps from 'matplotlib'*. Retrieved from
776 <https://CRAN.R-project.org/package=viridis>

- 777 Gelman, A., & Su, Y.-S. (2018). *Arm: Data analysis using regression and*
778 *multilevel/hierarchical models*. Retrieved from
779 <https://CRAN.R-project.org/package=arm>
- 780 Hair, J. F., Black, W. C., Babin, B. J., & Anderson, R. E. (2018). *Multivariate data analysis*.
781 Cengage.
- 782 Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data*
783 *mining, inference, and prediction* (Second Edition). New York, NY: Springer New
784 York.
- 785 Henry, L., & Wickham, H. (2019). *Purrr: Functional programming tools*. Retrieved from
786 <https://CRAN.R-project.org/package=purrr>
- 787 Higham, N. J. (2002). Computing the nearest correlation matrix—a problem from finance.
788 *IMA Journal of Numerical Analysis*, 22(3), 329–343.
789 <https://doi.org/10.1093/imanum/22.3.329>
- 790 Hong, S. (1999). Generating correlation matrices with model error for simulation studies in
791 factor analysis: A combination of the Tucker-Koopman-Linn model and Wijsman's
792 algorithm. *Behavior Research Methods, Instruments, & Computers*, 31(4), 727–730.
793 <https://doi.org/10.3758/BF03200754>
- 794 Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis.
795 *Psychometrika*, 30(2), 179–185. <https://doi.org/10.1007/BF02289447>
- 796 Jacqmin-Gadda, H., Sibillot, S., Proust, C., Molina, J.-M., & Thiébaut, R. (2007).
797 Robustness of the linear mixed model to misspecified error distribution.
798 *Computational Statistics & Data Analysis*, 51(10), 5142–5154.
799 <https://doi.org/10.1016/j.csda.2006.05.021>

- 800 Jamshidian, M., & Bentler, P. M. (1998). A quasi-Newton method for minimum trace factor
801 analysis. *Journal of Statistical Computation and Simulation*, 62(1-2), 73–89.
802 <https://doi.org/10.1080/00949659808811925>
- 803 Jennrich, R. I. (2002). A simple general method for oblique rotation. *Psychometrika*, 67(1),
804 7–19. <https://doi.org/10.1007/BF02294706>
- 805 Knol, D. L., & Berger, M. P. (1991). Empirical comparison between factor analysis and
806 multidimensional item response models. *Multivariate Behavioral Research*, 26(3),
807 457–477. https://doi.org/10.1207/s15327906mbr2603_5
- 808 Knowles, J. E., & Frederick, C. (2019). *MerTools: Tools for analyzing mixed effect regression
809 models*. Retrieved from <https://CRAN.R-project.org/package=merTools>
- 810 Koller, M. (2016). robustlmm: An R package for robust estimation of linear mixed-effects
811 models. *Journal of Statistical Software*, 75(6), 1–24.
812 <https://doi.org/10.18637/jss.v075.i06>
- 813 Kshirsagar, A. M. (1959). Bartlett decomposition and Wishart distribution. *The Annals of
814 Mathematical Statistics*, 30(1), 239–241. <https://doi.org/10.1214/aoms/1177706379>
- 815 Leifeld, P. (2013). texreg: Conversion of statistical model output in R to LaTeX and HTML
816 tables. *Journal of Statistical Software*, 55(8), 1–24.
817 <https://doi.org/10.18637/jss.v055.i08>
- 818 Li, Q., Li, D., & Qi, H. (2010). Newton's method for computing the nearest correlation
819 matrix with a simple upper bound. *Journal of Optimization Theory and Applications*,
820 147(3), 546–568. <https://doi.org/10.1007/s10957-010-9738-6>
- 821 Lorenzo-Seva, U., & Ferrando, P. J. (2020). Not positive definite correlation matrices in
822 exploratory item factor analysis: Causes, consequences and a proposed solution.

- 823 *Structural Equation Modeling: A Multidisciplinary Journal*, 1–10.
- 824 <https://doi.org/10.1080/10705511.2020.1735393>
- 825 Lurie, P. M., & Goldberg, M. S. (1998). An approximate method for sampling correlated
826 random variables from partially-specified distributions. *Management Science*, 44(2),
827 203–218. <https://doi.org/10.1287/mnsc.44.2.203>
- 828 MacCallum, R. C., & Tucker, L. R. (1991). Representing sources of error in the
829 common-factor model: Implications for theory and practice. *Psychological Bulletin*,
830 109(3), 502. <https://doi.org/10.1037/0033-2909.109.3.502>
- 831 MacCallum, R. C., Widaman, K. F., Preacher, K. J., & Hong, S. (2001). Sample size in
832 factor analysis: The role of model error. *Multivariate Behav. Res.*, 36(4), 611–637.
833 https://doi.org/10.1207/S15327906MBR3604_06
- 834 Maechler, M. (2019). *Sfsmisc: Utilities from 'seminar fuer statistik' eth zurich*. Retrieved
835 from <https://CRAN.R-project.org/package=sfsmisc>
- 836 Magnus, J. R., & Neudecker, H. (2019). *Matrix differential calculus with applications in
837 statistics and econometrics*. John Wiley & Sons.
838 <https://doi.org/10.1002/9781119541219>
- 839 Marwick, B. (2019). *Wordcountaddin: Word counts and readability statistics in r markdown
840 documents*. Retrieved from <https://github.com/benmarwick/wordcountaddin>
- 841 Meschiari, S. (2015). *Latex2exp: Use latexexpressions in plots*. Retrieved from
842 <https://CRAN.R-project.org/package=latex2exp>
- 843 Michalke, M. (2018a). *KoRpus: An r package for text analysis*. Retrieved from
844 <https://reaktanz.de/?c=hacking&s=koRpus>
- 845 Michalke, M. (2018b). *Syll: Hyphenation and syllable counting for text analysis*. Retrieved

- 846 from <https://reaktanz.de/?c=hacking&s=syll>
- 847 Michalke, M. (2019). *KoRpus.lang.en: Language support for 'koRpus' package: English.*
- 848 Retrieved from <https://reaktanz.de/?c=hacking&s=koRpus>
- 849 Mulaik, S. A. (2009). *Foundations of factor analysis*. Chapman; Hall/CRC.
- 850 <https://doi.org/10.1201/b15851>
- 851 Müller, K. (2017). *Here: A simpler way to find your files*. Retrieved from
- 852 <https://CRAN.R-project.org/package=here>
- 853 Müller, K., & Wickham, H. (2019). *Tibble: Simple data frames*. Retrieved from
- 854 <https://CRAN.R-project.org/package=tibble>
- 855 Olsson, U. (1979). Maximum likelihood estimation of the polychoric correlation coefficient.
- 856 *Psychometrika*, 44(4), 443–460. <https://doi.org/10.1007/BF02296207>
- 857 Pedersen, T. L. (2019). *Patchwork: The composer of plots*. Retrieved from
- 858 <https://CRAN.R-project.org/package=patchwork>
- 859 Qi, H., & Sun, D. (2006). A quadratically convergent Newton method for computing the
- 860 nearest correlation matrix. *SIAM Journal on Matrix Analysis and Applications*,
- 861 28(2), 360–385. <https://doi.org/10.1137/050624509>
- 862 R Core Team. (2019). *R: A language and environment for statistical computing*. Vienna,
- 863 Austria: R Foundation for Statistical Computing. Retrieved from
- 864 <https://www.R-project.org/>
- 865 Roff, M. (1936). Some properties of the communality in multiple factor theory.
- 866 *Psychometrika*, 1(2), 1–6. <https://doi.org/10.1007/BF02287999>
- 867 Sarkar, D. (2008). *Lattice: Multivariate data visualization with R*. New York: Springer.

- 868 <https://doi.org/10.1007/978-0-387-75969-2>
- 869 Thurstone, L. L. (1947). *Multiple-factor analysis; a development and expansion of the*
870 *vectors of mind* (pp. xix, 535). University of Chicago Press.
- 871 Tucker, L. R., Koopman, R. F., & Linn, R. L. (1969). Evaluation of factor analytic research
872 procedures by means of simulated correlation matrices. *Psychometrika*, 34(4),
873 421–459. <https://doi.org/10.1007/BF02290601>
- 874 Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* (Fourth). New
875 York: Springer. <https://doi.org/10.1007/978-0-387-21706-2>
- 876 Verbeke, G., & Lesaffre, E. (1997). The effect of misspecifying the random-effects
877 distribution in linear mixed models for longitudinal data. *Computational Statistics &*
878 *Data Analysis*, 23(4), 541–556. [https://doi.org/10.1016/S0167-9473\(96\)00047-3](https://doi.org/10.1016/S0167-9473(96)00047-3)
- 879 Waller, N. G. (2019). *Fungible: Psychometric functions from the Waller lab*.
- 880 Wickham, H. (2016). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
881 <https://doi.org/10.1007/978-3-319-24277-4>
- 882 Wickham, H. (2019a). *Forcats: Tools for working with categorical variables (factors)*.
883 Retrieved from <https://CRAN.R-project.org/package=forcats>
- 884 Wickham, H. (2019b). *Stringr: Simple, consistent wrappers for common string operations*.
885 Retrieved from <https://CRAN.R-project.org/package=stringr>
- 886 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ... Yutani,
887 H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686.
888 <https://doi.org/10.21105/joss.01686>
- 889 Wickham, H., François, R., Henry, L., & Müller, K. (2019). *Dplyr: A grammar of data*

- 890 *manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- 891 Wickham, H., & Henry, L. (2019). *Tidyr: Tidy messy data*. Retrieved from
892 <https://CRAN.R-project.org/package=tidyr>
- 893 Wickham, H., Hester, J., & Francois, R. (2018). *Readr: Read rectangular text data*.
894 Retrieved from <https://CRAN.R-project.org/package=readr>
- 895 Widaman, K. F., & Herringer, L. G. (1985). Iterative least squares estimates of
896 communality: Initial estimate need not affect stabilized value. *Psychometrika*, 50(4),
897 469–477. <https://doi.org/10.1007/BF02296264>
- 898 Wirth, R., & Edwards, M. C. (2007). Item factor analysis: Current approaches and future
899 directions. *Psychological Methods*, 12(1), 58.
900 <https://doi.org/10.1037/1082-989X.12.1.58>
- 901 Wotheke, W. (1993). Testing structural equation models. In K. A. Bollen & J. S. Long (Eds.),
902 *Testing structural equation models* (pp. 256–293). A Sage Focus Edition.
- 903 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida:
904 Chapman; Hall/CRC. <https://doi.org/10.1201/b15166>
- 905 Zhang, D., & Davidian, M. (2001). Linear mixed models with flexible distributions of
906 random effects for longitudinal data. *Biometrics*, 57(3), 795–802.
907 <https://doi.org/10.1111/j.0006-341X.2001.00795.x>

Table 1

Number of items (p) and subjects (N) resulting from each combination of number of factors (Factors), number of items per factor (Items/Factor), and subjects per item (Subjects/Item).

| Factors | Items/Factor | Subjects/Item | p | N |
|---------|--------------|---------------|-----|-------|
| 1 | 5 | 5 | 5 | 25 |
| 3 | 5 | 5 | 15 | 75 |
| 5 | 5 | 5 | 25 | 125 |
| 10 | 5 | 5 | 50 | 250 |
| 1 | 10 | 5 | 10 | 50 |
| 3 | 10 | 5 | 30 | 150 |
| 5 | 10 | 5 | 50 | 250 |
| 10 | 10 | 5 | 100 | 500 |
| 1 | 5 | 10 | 5 | 50 |
| 3 | 5 | 10 | 15 | 150 |
| 5 | 5 | 10 | 25 | 250 |
| 10 | 5 | 10 | 50 | 500 |
| 1 | 10 | 10 | 10 | 100 |
| 3 | 10 | 10 | 30 | 300 |
| 5 | 10 | 10 | 50 | 500 |
| 10 | 10 | 10 | 100 | 1,000 |
| 1 | 5 | 15 | 5 | 75 |
| 3 | 5 | 15 | 15 | 225 |
| 5 | 5 | 15 | 25 | 375 |
| 10 | 5 | 15 | 50 | 750 |
| 1 | 10 | 15 | 10 | 150 |
| 3 | 10 | 15 | 30 | 450 |
| 5 | 10 | 15 | 50 | 750 |
| 10 | 10 | 15 | 100 | 1,500 |

Table 2

Coefficient estimates and standard errors for the linear and polynomial mixed effects models using $\log[D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})]$ as the dependent variable and estimating a random intercept for each indefinite correlation matrix.

| | Linear Model | Polynomial Model |
|---|----------------|------------------|
| Constant | -2.209 (0.001) | -2.339 (0.001) |
| Subjects/Item | -0.300 (0.001) | -0.249 (0.001) |
| Items/Factor | -0.229 (0.001) | -0.262 (0.001) |
| Factors | -0.371 (0.001) | -0.521 (0.001) |
| Factor Loading | -0.048 (0.001) | -0.048 (0.001) |
| Model Error | -0.008 (0.001) | -0.010 (0.001) |
| Smoothing Method (APA) | -0.015 (0.000) | -0.009 (0.000) |
| Smoothing Method (BY) | -0.067 (0.000) | -0.058 (0.000) |
| Smoothing Method (KB) | -0.020 (0.000) | -0.013 (0.000) |
| Subjects/Item ² | | 0.078 (0.002) |
| Factors ² | | 0.189 (0.001) |
| Model Error ² | | -0.004 (0.002) |
| Subjects/Item × Items/Factor | -0.006 (0.001) | -0.005 (0.001) |
| Subjects/Item × Factors | 0.016 (0.001) | 0.005 (0.001) |
| Subjects/Item × Factor Loading | 0.007 (0.001) | -0.027 (0.001) |
| Subjects/Item × Model Error | -0.001 (0.001) | -0.004 (0.001) |
| Subjects/Item × Smoothing Method (APA) | 0.019 (0.000) | 0.017 (0.000) |
| Subjects/Item × Smoothing Method (BY) | 0.038 (0.000) | 0.032 (0.000) |
| Subjects/Item × Smoothing Method (KB) | 0.027 (0.000) | 0.025 (0.000) |
| Subjects/Item × Factors ² | | -0.006 (0.001) |
| Subjects/Item × Model Error ² | | -0.000 (0.001) |
| Items/Factor × Factors | -0.034 (0.001) | 0.009 (0.001) |
| Items/Factor × Factor Loading | -0.005 (0.001) | 0.001 (0.001) |
| Items/Factor × Model Error | 0.002 (0.001) | 0.001 (0.001) |
| Items/Factor × Smoothing Method (APA) | -0.000 (0.000) | -0.000 (0.000) |
| Items/Factor × Smoothing Method (BY) | 0.018 (0.000) | 0.016 (0.000) |
| Items/Factor × Smoothing Method (KB) | 0.000 (0.000) | -0.000 (0.000) |
| Items/Factor × Subjects/Item ² | | 0.003 (0.001) |
| Items/Factor × Factors ² | | -0.003 (0.001) |
| Items/Factor × Model Error ² | | -0.000 (0.001) |
| Factors × Factor Loading | 0.033 (0.001) | 0.077 (0.001) |
| Factors × Model Error | 0.001 (0.001) | -0.002 (0.001) |
| Factors × Smoothing Method (APA) | 0.002 (0.000) | 0.003 (0.000) |
| Factors × Smoothing Method (BY) | -0.005 (0.000) | -0.014 (0.000) |
| Factors × Smoothing Method (KB) | -0.000 (0.000) | -0.002 (0.000) |
| Factors × Subjects/Item ² | | 0.002 (0.001) |
| Factors × Model Error ² | | -0.001 (0.001) |

| | Linear Model | Polynomial Model |
|--|----------------|------------------|
| Factor Loading \times Model Error | 0.009 (0.001) | 0.008 (0.001) |
| Factor Loading \times Smoothing Method (APA) | -0.008 (0.000) | -0.008 (0.000) |
| Factor Loading \times Smoothing Method (BY) | 0.024 (0.000) | 0.024 (0.000) |
| Factor Loading \times Smoothing Method (KB) | -0.011 (0.000) | -0.011 (0.000) |
| Factor Loading \times Subjects/Item ² | | 0.002 (0.001) |
| Factor Loading \times Factors ² | | -0.049 (0.001) |
| Factor Loading \times Model Error ² | | 0.003 (0.001) |
| Model Error \times Smoothing Method (APA) | -0.003 (0.000) | -0.003 (0.000) |
| Model Error \times Smoothing Method (BY) | 0.001 (0.000) | 0.000 (0.000) |
| Model Error \times Smoothing Method (KB) | -0.004 (0.000) | -0.004 (0.000) |
| Model Error \times Subjects/Item ² | | 0.001 (0.001) |
| Model Error \times Factors ² | | 0.002 (0.001) |
| Smoothing Method (APA) \times Subjects/Item ² | | -0.009 (0.000) |
| Smoothing Method (BY) \times Subjects/Item ² | | -0.028 (0.000) |
| Smoothing Method (KB) \times Subjects/Item ² | | -0.013 (0.000) |
| Smoothing Method (APA) \times Factors ² | | -0.002 (0.000) |
| Smoothing Method (BY) \times Factors ² | | 0.012 (0.000) |
| Smoothing Method (KB) \times Factors ² | | 0.002 (0.000) |
| Smoothing Method (APA) \times Model Error ² | | -0.001 (0.000) |
| Smoothing Method (BY) \times Model Error ² | | 0.001 (0.000) |
| Smoothing Method (KB) \times Model Error ² | | -0.001 (0.000) |
| Subjects/Item ² \times Factors ² | | 0.000 (0.001) |
| Subjects/Item ² \times Model Error ² | | 0.002 (0.002) |
| Factors ² \times Model Error ² | | 0.001 (0.001) |
| AIC | -1808591.524 | -1938579.038 |
| BIC | -1808191.308 | -1937878.660 |
| Log Likelihood | 904331.762 | 969352.519 |
| Num. obs. | 497381 | 497381 |
| Num. groups: id | 124346 | 124346 |
| Var: id (Intercept) | 0.017 | 0.008 |
| Var: Residual | 0.000 | 0.000 |

Table 3

Coefficient estimates and standard errors for the linear and polynomial mixed effects models using $\log[RMSE(\mathbf{F}, \hat{\mathbf{F}})]$ as the dependent variable and estimating a random intercept for each indefinite correlation matrix.

| | Linear Model | Polynomial Model |
|---|----------------|------------------|
| Constant | -2.235 (0.001) | -2.229 (0.003) |
| Subjects/Item | -0.189 (0.001) | -0.183 (0.003) |
| Items/Factor | -0.175 (0.001) | -0.369 (0.002) |
| Factors | -0.255 (0.001) | -0.344 (0.002) |
| Factor Loading | -0.438 (0.001) | -0.574 (0.002) |
| Model Error | 0.104 (0.001) | 0.229 (0.002) |
| Smoothing Method (APA) | -0.006 (0.000) | -0.004 (0.001) |
| Smoothing Method (BY) | -0.019 (0.000) | -0.033 (0.001) |
| Smoothing Method (KB) | -0.011 (0.000) | -0.008 (0.001) |
| Extraction Method (ML) | 0.110 (0.000) | -0.049 (0.001) |
| Extraction Method (PA) | 0.005 (0.000) | 0.004 (0.001) |
| Subjects/Item ² | | 0.157 (0.003) |
| Factor Loading ² | | 0.006 (0.003) |
| Factors ² | | 0.095 (0.002) |
| Model Error ² | | 0.099 (0.003) |
| Subjects/Item × Items/Factor | 0.016 (0.001) | 0.024 (0.001) |
| Subjects/Item × Factors | 0.021 (0.001) | 0.033 (0.001) |
| Subjects/Item × Factor Loading | -0.017 (0.001) | -0.076 (0.003) |
| Subjects/Item × Model Error | 0.025 (0.001) | 0.019 (0.001) |
| Subjects/Item × Smoothing Method (APA) | 0.003 (0.000) | 0.003 (0.000) |
| Subjects/Item × Smoothing Method (BY) | 0.003 (0.000) | 0.001 (0.000) |
| Subjects/Item × Smoothing Method (KB) | 0.005 (0.000) | 0.006 (0.000) |
| Subjects/Item × Extraction Method (ML) | 0.107 (0.000) | 0.146 (0.000) |
| Subjects/Item × Extraction Method (PA) | -0.000 (0.000) | -0.000 (0.000) |
| Subjects/Item × Factor Loading ² | | -0.004 (0.004) |
| Subjects/Item × Factors ² | | -0.020 (0.001) |
| Subjects/Item × Model Error ² | | 0.012 (0.002) |
| Items/Factor × Factors | -0.005 (0.001) | 0.041 (0.001) |
| Items/Factor × Factor Loading | -0.004 (0.001) | -0.052 (0.001) |
| Items/Factor × Model Error | 0.036 (0.001) | 0.054 (0.001) |
| Items/Factor × Smoothing Method (APA) | 0.002 (0.000) | 0.003 (0.000) |
| Items/Factor × Smoothing Method (BY) | 0.012 (0.000) | 0.017 (0.000) |
| Items/Factor × Smoothing Method (KB) | 0.002 (0.000) | 0.002 (0.000) |
| Items/Factor × Extraction Method (ML) | 0.082 (0.000) | 0.102 (0.000) |
| Items/Factor × Extraction Method (PA) | -0.003 (0.000) | -0.005 (0.000) |
| Items/Factor × Subjects/Item ² | | -0.005 (0.001) |
| Items/Factor × Factor Loading ² | | 0.144 (0.001) |

| | Linear Model | Polynomial Model |
|---|----------------|------------------|
| Items/Factor \times Factors ² | | -0.012 (0.001) |
| Items/Factor \times Model Error ² | | 0.021 (0.001) |
| Factors \times Factor Loading | -0.014 (0.001) | -0.014 (0.001) |
| Factors \times Model Error | 0.042 (0.001) | 0.068 (0.001) |
| Factors \times Smoothing Method (APA) | 0.001 (0.000) | 0.001 (0.000) |
| Factors \times Smoothing Method (BY) | 0.000 (0.000) | -0.003 (0.000) |
| Factors \times Smoothing Method (KB) | 0.001 (0.000) | -0.000 (0.000) |
| Factors \times Extraction Method (ML) | -0.090 (0.000) | -0.203 (0.000) |
| Factors \times Extraction Method (PA) | -0.004 (0.000) | -0.006 (0.000) |
| Factors \times Subjects/Item ² | | -0.007 (0.002) |
| Factors \times Factor Loading ² | | -0.073 (0.002) |
| Factors \times Model Error ² | | 0.026 (0.002) |
| Factor Loading \times Model Error | -0.047 (0.001) | -0.023 (0.001) |
| Factor Loading \times Smoothing Method (APA) | 0.000 (0.000) | 0.000 (0.000) |
| Factor Loading \times Smoothing Method (BY) | 0.021 (0.000) | 0.022 (0.000) |
| Factor Loading \times Smoothing Method (KB) | -0.001 (0.000) | -0.001 (0.000) |
| Factor Loading \times Extraction Method (ML) | 0.204 (0.000) | 0.217 (0.000) |
| Factor Loading \times Extraction Method (PA) | -0.003 (0.000) | -0.004 (0.000) |
| Factor Loading \times Subjects/Item ² | | 0.004 (0.003) |
| Factor Loading \times Factors ² | | 0.018 (0.001) |
| Factor Loading \times Model Error ² | | 0.003 (0.002) |
| Model Error \times Smoothing Method (APA) | 0.000 (0.000) | 0.000 (0.000) |
| Model Error \times Smoothing Method (BY) | 0.001 (0.000) | 0.001 (0.000) |
| Model Error \times Smoothing Method (KB) | -0.000 (0.000) | -0.000 (0.000) |
| Model Error \times Extraction Method (ML) | -0.034 (0.000) | -0.035 (0.000) |
| Model Error \times Extraction Method (PA) | -0.001 (0.000) | -0.001 (0.000) |
| Model Error \times Subjects/Item ² | | -0.034 (0.001) |
| Model Error \times Factor Loading ² | | -0.117 (0.001) |
| Model Error \times Factors ² | | -0.022 (0.001) |
| Smoothing Method (APA) \times Extraction Method (ML) | 0.006 (0.001) | 0.006 (0.000) |
| Smoothing Method (BY) \times Extraction Method (ML) | 0.019 (0.001) | 0.019 (0.000) |
| Smoothing Method (KB) \times Extraction Method (ML) | 0.011 (0.001) | 0.011 (0.000) |
| Smoothing Method (APA) \times Extraction Method (PA) | -0.002 (0.001) | -0.002 (0.000) |
| Smoothing Method (BY) \times Extraction Method (PA) | -0.004 (0.001) | -0.003 (0.000) |
| Smoothing Method (KB) \times Extraction Method (PA) | -0.002 (0.001) | -0.002 (0.000) |
| Smoothing Method (APA) \times Subjects/Item ² | | -0.002 (0.000) |
| Smoothing Method (BY) \times Subjects/Item ² | | -0.011 (0.000) |
| Smoothing Method (KB) \times Subjects/Item ² | | -0.004 (0.000) |
| Smoothing Method (APA) \times Factor Loading ² | | 0.001 (0.000) |
| Smoothing Method (BY) \times Factor Loading ² | | 0.016 (0.000) |
| Smoothing Method (KB) \times Factor Loading ² | | 0.001 (0.000) |
| Smoothing Method (APA) \times Factors ² | | 0.000 (0.000) |
| Smoothing Method (BY) \times Factors ² | | 0.005 (0.000) |
| Smoothing Method (KB) \times Factors ² | | 0.001 (0.000) |

| | Linear Model | Polynomial Model |
|---|--------------|------------------|
| Smoothing Method (APA) \times Model Error ² | | -0.000 (0.000) |
| Smoothing Method (BY) \times Model Error ² | | -0.001 (0.000) |
| Smoothing Method (KB) \times Model Error ² | | -0.000 (0.000) |
| Extraction Method (ML) \times Subjects/Item ² | | 0.011 (0.000) |
| Extraction Method (PA) \times Subjects/Item ² | | 0.001 (0.000) |
| Extraction Method (ML) \times Factor Loading ² | | 0.105 (0.000) |
| Extraction Method (PA) \times Factor Loading ² | | 0.001 (0.000) |
| Extraction Method (ML) \times Factors ² | | 0.139 (0.000) |
| Extraction Method (PA) \times Factors ² | | 0.004 (0.000) |
| Extraction Method (ML) \times Model Error ² | | -0.016 (0.000) |
| Extraction Method (PA) \times Model Error ² | | -0.000 (0.000) |
| Subjects/Item ² \times Factor Loading ² | | -0.087 (0.004) |
| Subjects/Item ² \times Factors ² | | 0.010 (0.002) |
| Subjects/Item ² \times Model Error ² | | -0.003 (0.002) |
| Factor Loading ² \times Factors ² | | 0.040 (0.002) |
| Factor Loading ² \times Model Error ² | | -0.052 (0.002) |
| Factors ² \times Model Error ² | | -0.027 (0.002) |
| AIC | -2414463.669 | -2801683.227 |
| BIC | -2413804.118 | -2800461.837 |
| Log Likelihood | 1207285.834 | 1400941.614 |
| Num. obs. | 1489425 | 1489425 |
| Num. groups: id | 124346 | 124346 |
| Var: id (Intercept) | 0.032 | 0.017 |
| Var: Residual | 0.008 | 0.007 |

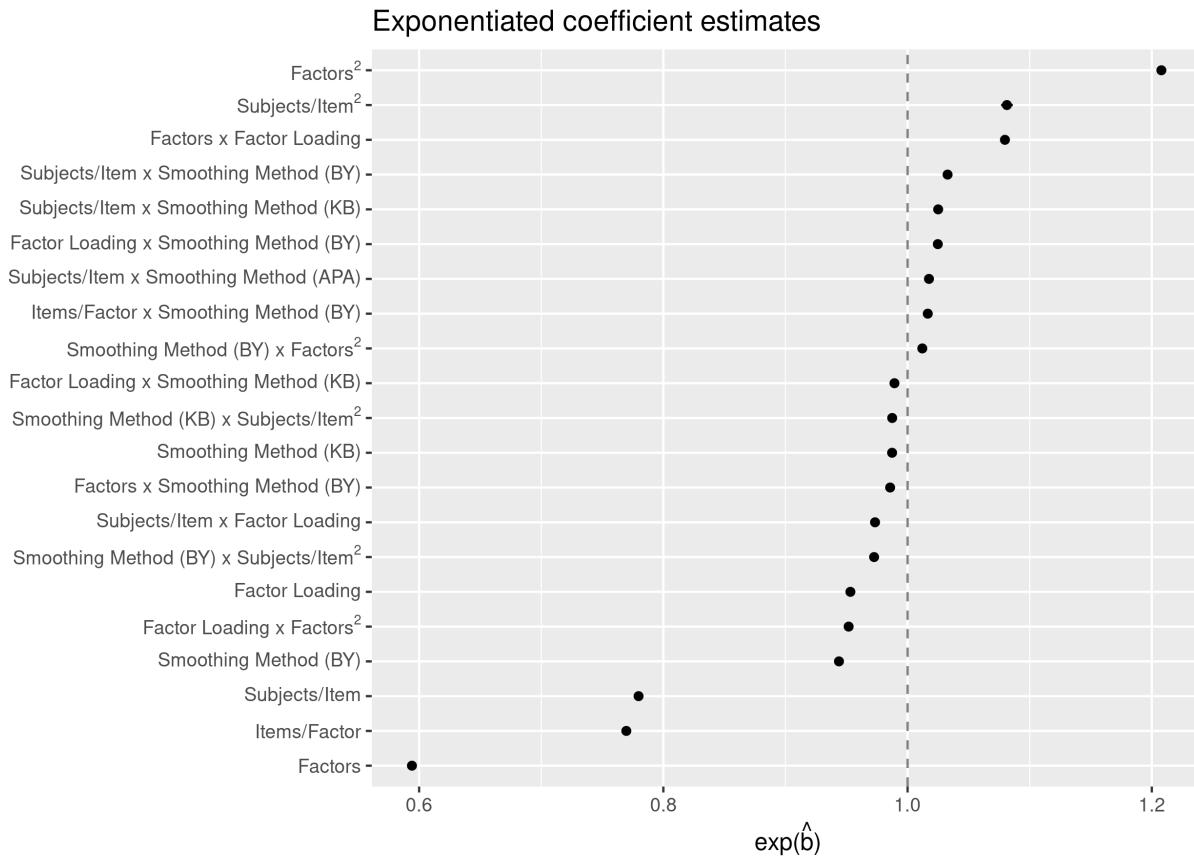


Figure 1. Exponentiated coefficient estimates for the mixed effects model using $\log[D_s(\mathbf{R}_{Sm}, \mathbf{R}_{Pop})]$ as the dependent variable (Model 1B). APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991). The effect of the condition where no smoothing was applied is subsumed within the Constant term.

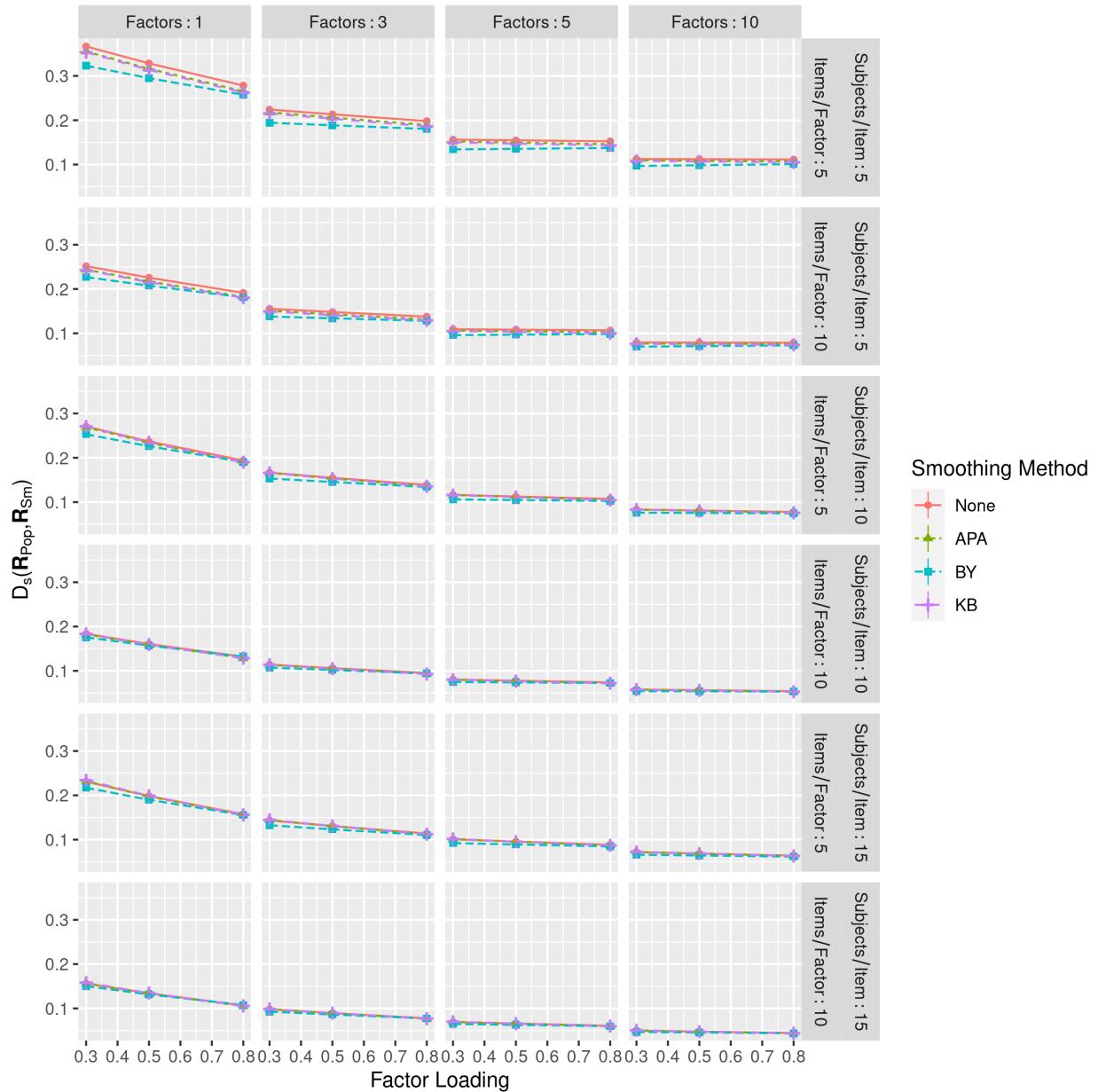


Figure 2. Scaled distance between the smoothed (\mathbf{R}_{Sm}) and model-implied (\mathbf{R}_{Pop}) correlation matrices. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing.

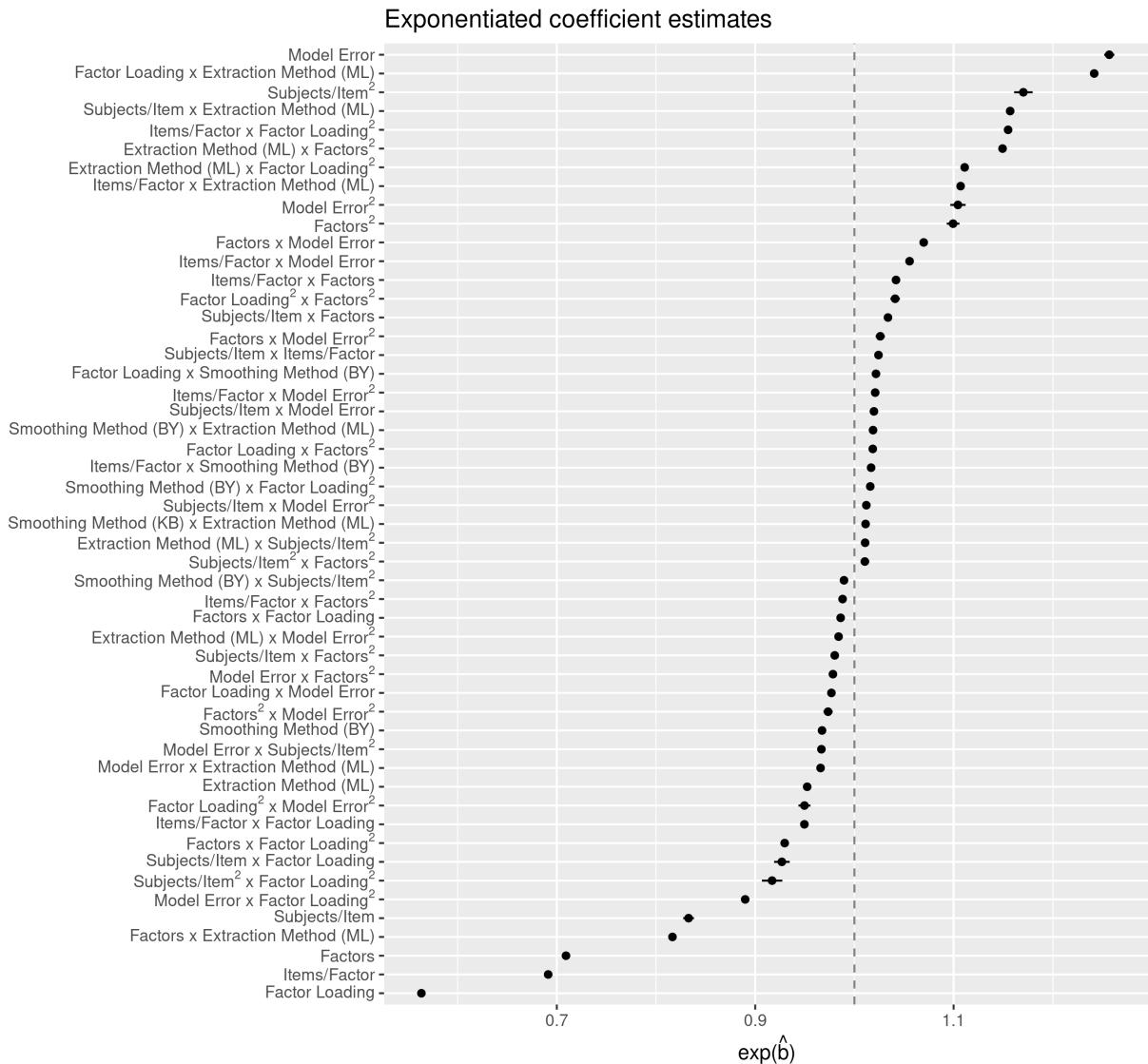


Figure 3. Exponentiated coefficient estimates for the mixed effects model using $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$ as the dependent variable (Model 2B). APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); ML = Maximum likelihood; PA = Principal axis. The effects of no smoothing and ordinary least squares factor analysis are subsumed within the Constant term.

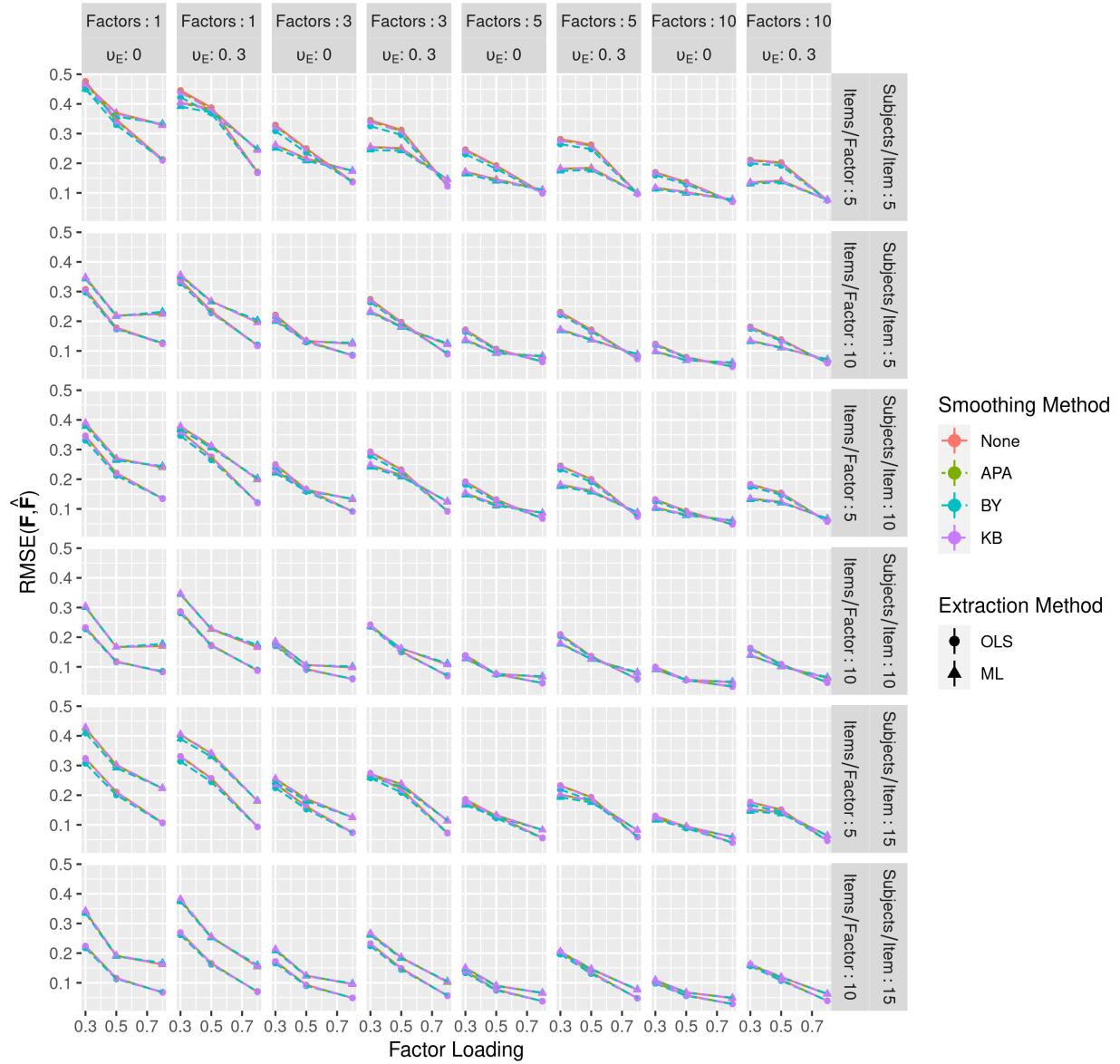


Figure 4. Estimated marginal mean $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ values and 99% confidence intervals. To conserve space, the intermediate values of model error and subjects per item have been omitted. The principal axis factor extraction method was also omitted because it led to nearly identical results compared to ordinary least squares. OLS = ordinary least squares; ML = maximum likelihood; APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991).

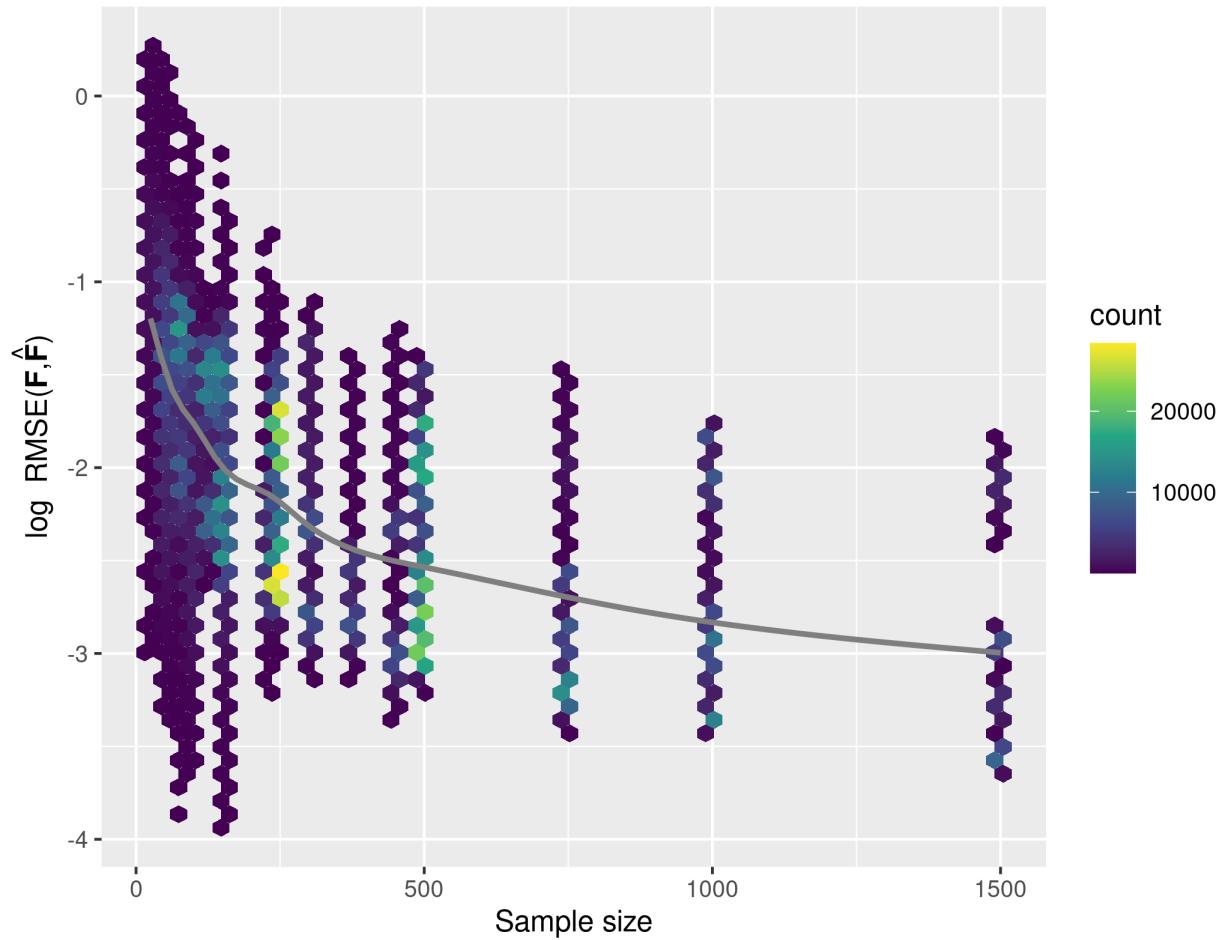


Figure 5. Log root-mean-square error (RMSE) between the true and estimated factor loading matrices as a function of sample size.

Appendix A

Regression Diagnostics

908 Models 1A and 1B: Regression models predicting $\log D_s(\mathbf{R}_{\text{Pop}}, \mathbf{R}_{\text{Sm}})$

909 Models 1A and 1B were linear mixed-effects models predicting the (log) scaled
 910 distance between the smoothed and model-implied population correlation matrix and was fit
 911 using the R *lme4* package (Version 1.1.23; Bates, Mächler, Bolker, & Walker, 2015). Model
 912 1A was a linear model fit using all simulation variables and their interactions. In Model 1B,
 913 second-degree polynomial terms were added for number of factors, number of subjects per
 914 item, factor loading, and model error. Diagnostic plots showing standardized residuals
 915 plotted against fitted values for both models, quantile-quantile (QQ) plots of the residuals,
 916 and QQ plots for the random intercept terms are shown in Figures A3, A1, and A2
 917 respectively. These plots show that some assumptions of the linear mixed-effects model seem
 918 to have been violated for Models 1A and 1B, even after applying a log-transformation to the
 919 response variable.

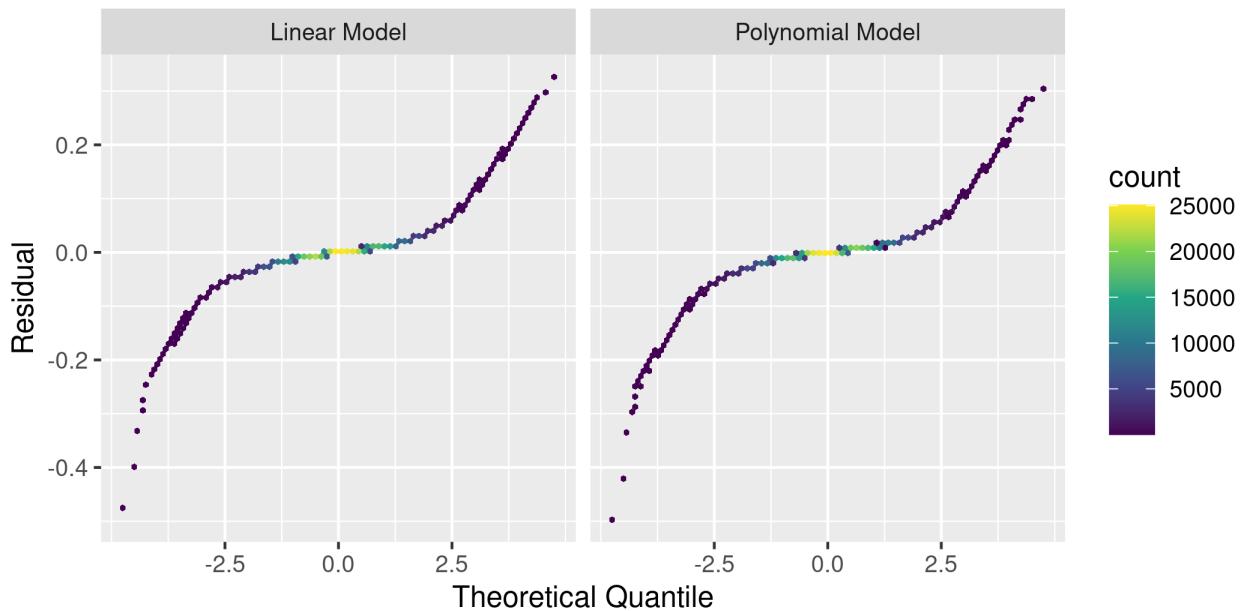


Figure A1. Quantile-quantile plot of residuals for Models 1A and 1B.

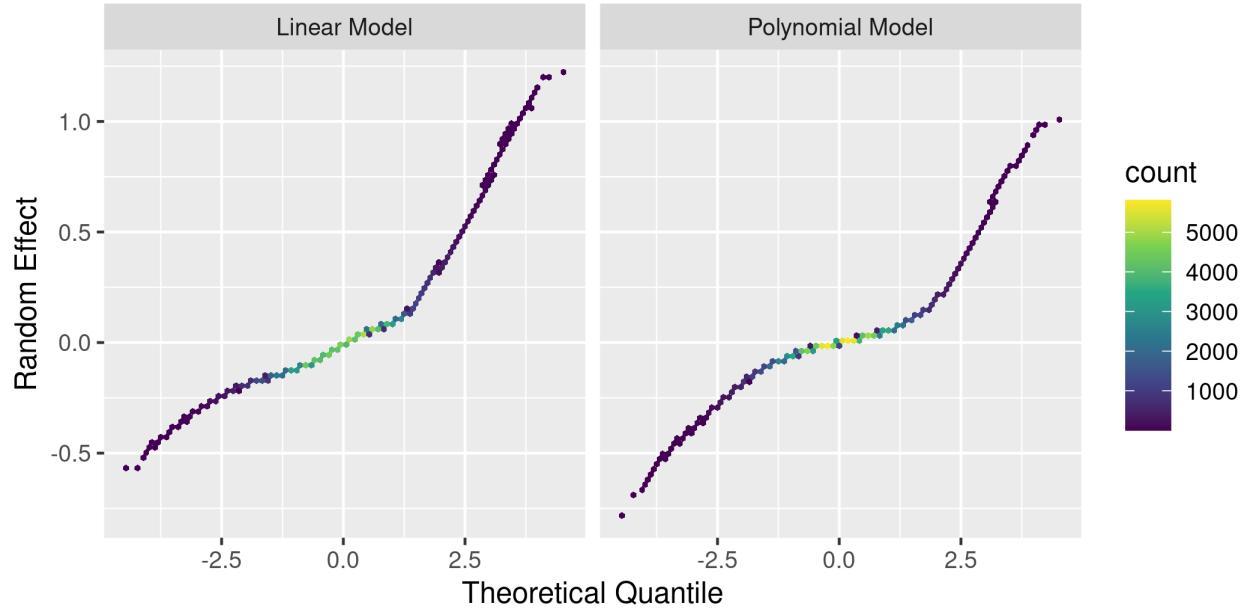


Figure A2. Quantile-quantile plot of random intercept terms for Models 1A and 1B.

920 Figure A3 shows that the variance of the residuals was not constant over the range of
 921 fitted values for both the linear and polynomial models. In particular, for both models there
 922 was little variation near the edges of the range of fitted values and a large amount of
 923 variation near the center of the distribution of fitted values. Therefore, the homoscedasticity
 924 assumption seemed to have been violated. Moreover, Figure A1 shows that the assumption
 925 of normally-distributed errors was also likely violated. In particular, Figure A1 shows that
 926 the distributions of residuals (for both models) had heavy tails and had a slight positive
 927 skew (Model 1A: kurtosis = 16.25, skew = 0.60; Model 1B: kurtosis = 18.61, skew = 0.23).
 928 Finally, Figure A2 shows that the random effects (random intercepts) were not
 929 normally-distributed for either the linear or polynomial model (Model 1A: kurtosis = 5.52,
 930 skew = 1.52; Model 1B: kurtosis = 10.33, skew = 0.59). To address these violations of the
 931 model assumptions, I first attempted to fit a robust mixed-effects model using `rlmer()`
 932 function in the R *robustlmm* package (Version 2.3; Koller, 2016). Unfortunately, the data set
 933 was too large for the `rlmer()` function to handle. I also tried a more complex
 934 transformation of the dependent variable (using a Box-Cox power transformation; Box &

935 Cox, 1964), but it produced no discernible benefit compared to a log transformation.

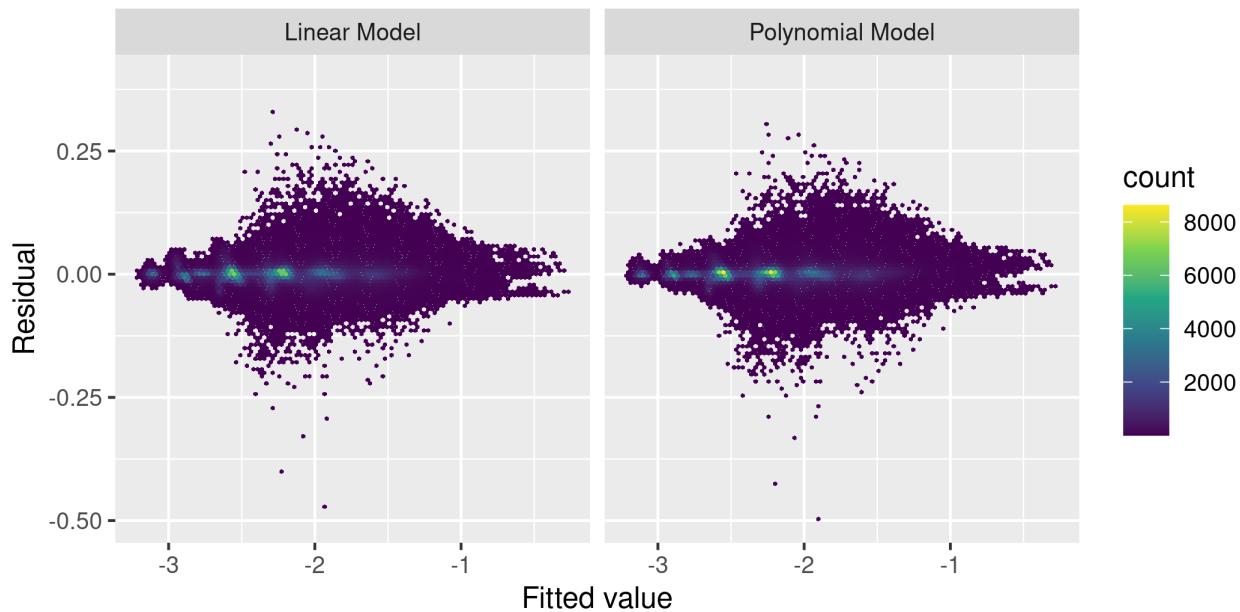


Figure A3. Residuals plotted against fitted values for Models 1A and 1B.

936 The apparent violations of the assumptions of the mixed-effects model were concerning.
 937 However, inference for the fixed effects in mixed-effects models seems to be somewhat robust
 938 to these violations. In particular, Jacqmin-Gadda, Sibillot, Proust, Molina, & Thiébaut
 939 (2007) showed that inference for fixed effects is robust for non-Gaussian and heteroscedastic
 940 errors. Moreover, Jacqmin-Gadda et al. (2007) cited several studies indicating that inference
 941 for fixed effects is also robust to non-Gaussian random effects (Butler & Louis, 1992; Verbeke
 942 & Lesaffre, 1997; Zhang & Davidian, 2001). Finally, the purpose of the present analysis was
 943 to obtain estimates of the fixed effects of matrix smoothing methods (and the interactions
 944 between smoothing methods and the other design factors) on population correlation matrix
 945 recovery. Neither p -values nor confidence intervals were of primary concern. Therefore, the
 946 apparent violation of some model assumptions likely did not affect the main results of this
 947 study.

948 **Models 2A and 2B: Regression models predicting $\log \text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$**

949 Models 2A and 2B were mixed-effects models predicting $\log \text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ and fit using
 950 the R *lme4* package (Bates, Mächler, Bolker, & Walker, 2015). Model 2A was a linear model
 951 fit using all simulation variables and their interactions. In Model 2B, second-degree
 952 polynomial terms were added for number of factors, number of subjects per item, factor
 953 loading, and model error. As with Models 1A and 1B, diagnostic plots showing standardized
 954 residuals plotted against fitted values for both models, QQ plots for the residuals, and QQ
 955 plots for the random intercept terms are shown in Figures A3, A5, and A6 respectively.

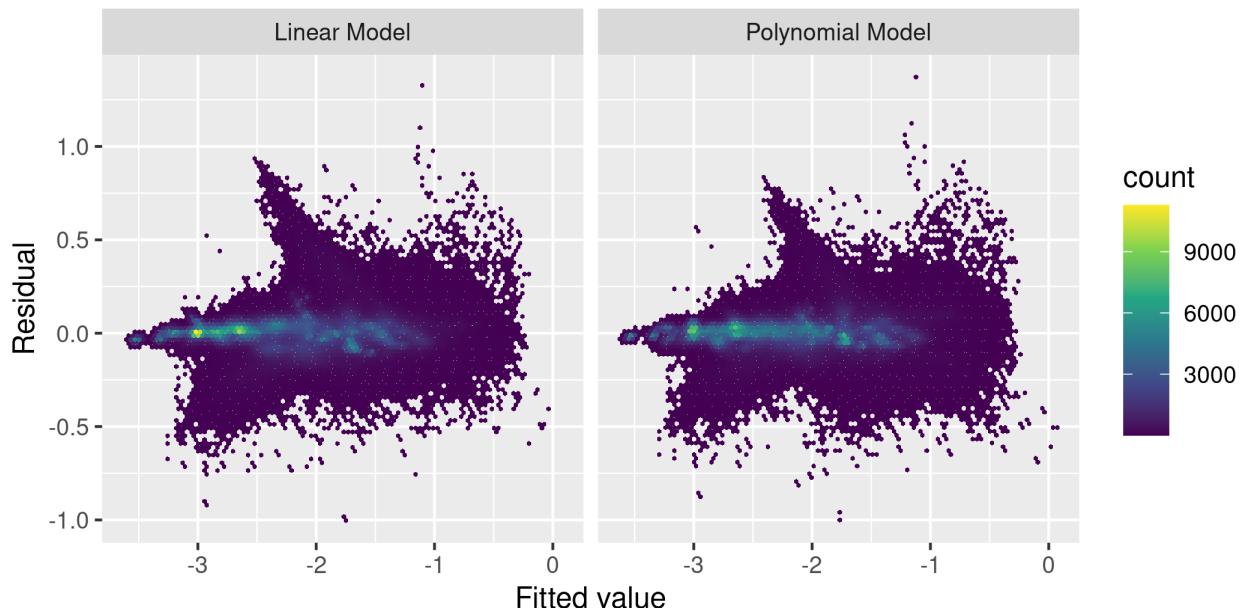


Figure A4. Residuals plotted against fitted values for Models 2A and 2B.

956 These plots indicate many of the same issues in Models 2A and 2B as were seen for
 957 Models 1A and 1B. First, Figure A3 shows clear evidence of non-homogeneous conditional
 958 error variance for both the linear and polynomial models. Specifically, the residual variance
 959 seemed generally to be larger for larger fitted values. Second, Figure A5 shows that the
 960 distribution of residuals for both models was non-normal and similar to the distributions of
 961 the residuals from Model 1A and 1B (i.e., positively-skewed and having heavy tails). Finally,
 962 Figure A6 shows that the estimated random effects were likewise not normally-distributed.

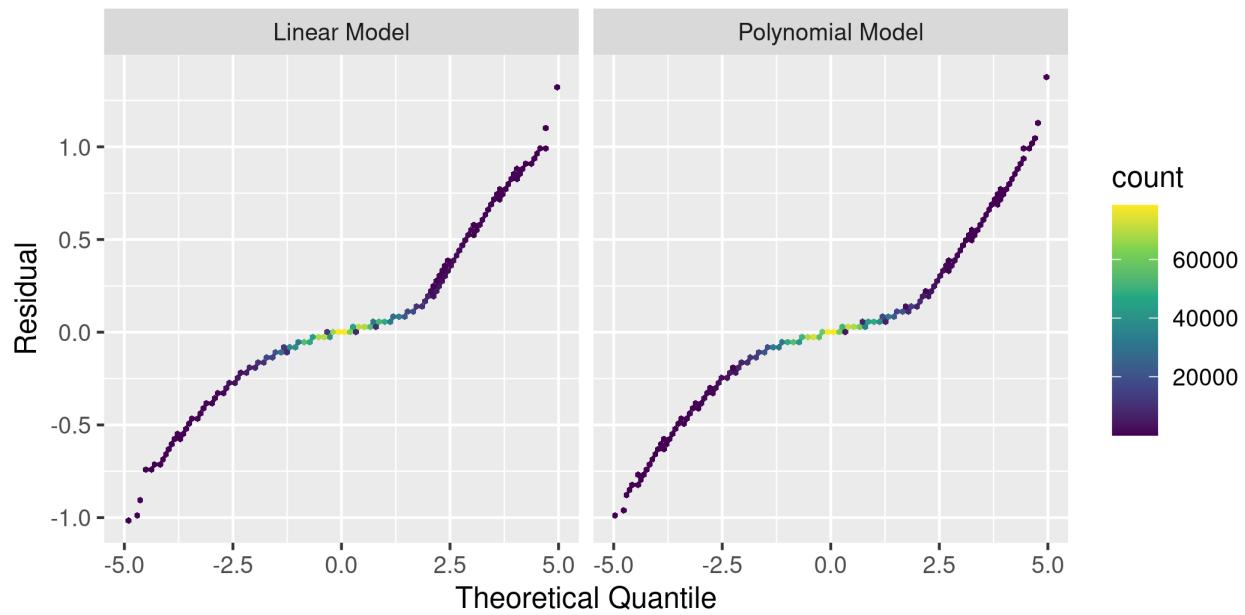


Figure A5. Quantile-quantile plot of residuals for Models 2A and 2B.

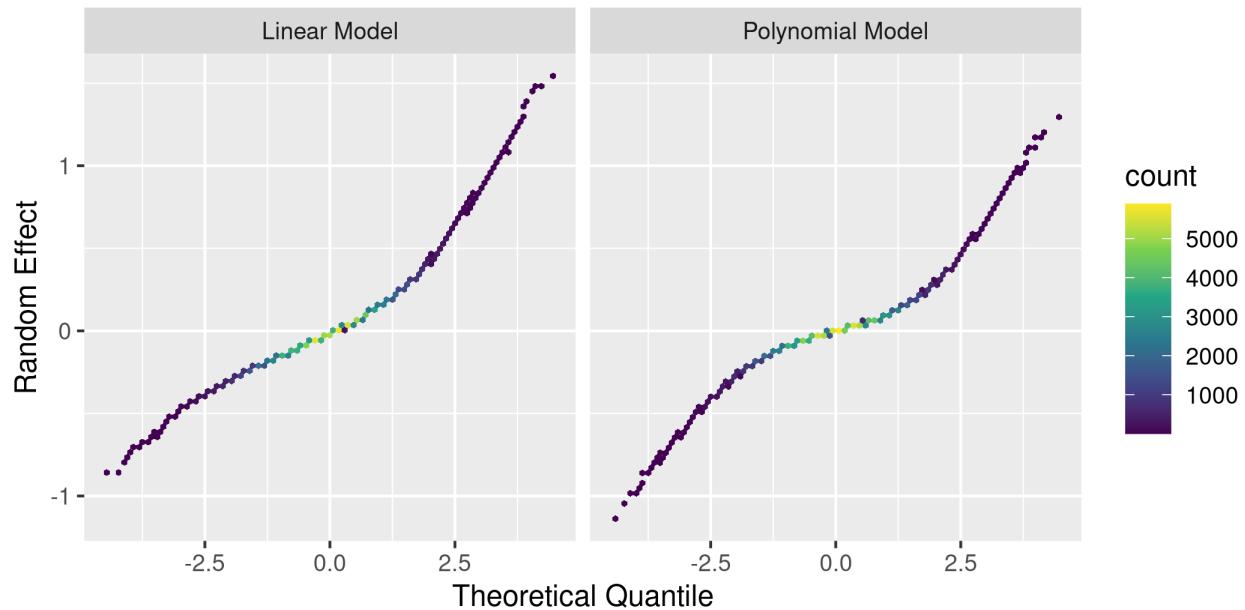


Figure A6. Quantile-quantile plot of random intercept terms for Models 2A and 2B.

963 The distribution of random intercepts was positively-skewed with heavy tails. Alternative
964 transformations of the dependent variable were tried but did not seem to improve model fit
965 compared to a log transformation. As with Model 1, these violations of the model
966 assumptions are somewhat concerning and indicate that the estimated parameters—the
967 estimated standard errors, in particular—should be treated with some degree of skepticism.
968 However, the main results of the study are unlikely to have been affected greatly by these
969 violations of the model assumptions.

Appendix B
Supplemental Tables and Figures

970 Indefinite Matrix Frequency

971 The percent of indefinite tetrachoric correlation matrices differed from condition to
972 condition. Table B1 reports the percent of indefinite matrices for each of the 216 conditions
973 of the study design. One of the more obvious trends in this table is that conditions with
974 more (major) factors tended to produce more indefinite tetrachoric correlation matrices.
975 Based on the results reported by Debelak and Tran (2013; 2016), who found that indefinite
976 tetrachoric and polychoric correlation matrices were much more common for data sets with
977 many items, this is likely due to the correlation of factor number with total number of items.
978 (See Lorenzo-Seva and Ferrando, 2020, for further discussion of the relationship between the
979 number of items and matrix indefiniteness.) Moreover, the results in Table B1 indicate that
980 indefinite matrices were more common for conditions with more items per factor, fewer
981 subjects per item, and higher factor loadings. All of these trends corroborate the similar
982 results by Debelak and Tran (2013; 2016) and their conclusions about which variables most
983 affected the frequency of indefinite tetrachoric or polychoric correlation matrices.

Table B1

Percent of indefinite tetrachoric correlation matrices by Number of Subjects Per Item (N/p), Number of Items per Factor (p/m), Factor Loading, Model Error (v_E), and Number of Factors.

| N/p | p/m | Loading | v_E | Factors | | | |
|-------|-------|---------|-------|---------|-------|-------|-------|
| | | | | 1 | 3 | 5 | 10 |
| 5 | 5 | 0.3 | 0.0 | 10.5 | 96.6 | 100.0 | 100.0 |
| 5 | 5 | 0.3 | 0.1 | 10.6 | 97.3 | 100.0 | 100.0 |
| 5 | 5 | 0.3 | 0.3 | 13.5 | 99.3 | 100.0 | 100.0 |
| 5 | 5 | 0.5 | 0.0 | 15.6 | 98.9 | 100.0 | 100.0 |
| 5 | 5 | 0.5 | 0.1 | 14.4 | 99.0 | 100.0 | 100.0 |
| 5 | 5 | 0.5 | 0.3 | 15.6 | 100.0 | 100.0 | 100.0 |
| 5 | 5 | 0.8 | 0.0 | 13.5 | 100.0 | 100.0 | 100.0 |
| 5 | 5 | 0.8 | 0.1 | 13.4 | 100.0 | 100.0 | 100.0 |
| 5 | 5 | 0.8 | 0.3 | 13.9 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.3 | 0.0 | 78.0 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.3 | 0.1 | 79.1 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.3 | 0.3 | 85.5 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.5 | 0.0 | 88.1 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.5 | 0.1 | 89.3 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.5 | 0.3 | 94.1 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.8 | 0.0 | 98.9 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.8 | 0.1 | 99.3 | 100.0 | 100.0 | 100.0 |
| 5 | 10 | 0.8 | 0.3 | 99.6 | 100.0 | 100.0 | 100.0 |
| 10 | 5 | 0.3 | 0.0 | 2.5 | 7.9 | 9.4 | 5.8 |
| 10 | 5 | 0.3 | 0.1 | 2.0 | 10.5 | 12.3 | 18.1 |
| 10 | 5 | 0.3 | 0.3 | 3.3 | 26.8 | 54.6 | 93.1 |
| 10 | 5 | 0.5 | 0.0 | 3.4 | 21.3 | 29.0 | 49.1 |
| 10 | 5 | 0.5 | 0.1 | 3.5 | 26.0 | 38.1 | 70.8 |
| 10 | 5 | 0.5 | 0.3 | 4.4 | 48.8 | 82.5 | 99.7 |
| 10 | 5 | 0.8 | 0.0 | 13.1 | 98.4 | 100.0 | 100.0 |
| 10 | 5 | 0.8 | 0.1 | 11.4 | 98.2 | 100.0 | 100.0 |
| 10 | 5 | 0.8 | 0.3 | 12.5 | 99.3 | 100.0 | 100.0 |
| 10 | 10 | 0.3 | 0.0 | 8.7 | 8.0 | 7.9 | 5.1 |
| 10 | 10 | 0.3 | 0.1 | 11.0 | 14.0 | 19.5 | 38.2 |
| 10 | 10 | 0.3 | 0.3 | 21.2 | 70.2 | 94.4 | 100.0 |
| 10 | 10 | 0.5 | 0.0 | 23.8 | 39.5 | 63.3 | 94.8 |
| 10 | 10 | 0.5 | 0.1 | 24.2 | 56.3 | 83.7 | 99.9 |
| 10 | 10 | 0.5 | 0.3 | 39.8 | 94.4 | 100.0 | 100.0 |
| 10 | 10 | 0.8 | 0.0 | 84.2 | 100.0 | 100.0 | 100.0 |
| 10 | 10 | 0.8 | 0.1 | 83.3 | 100.0 | 100.0 | 100.0 |

| | | | | | | | |
|----|----|-----|-----|------|-------|-------|-------|
| 10 | 10 | 0.8 | 0.3 | 89.9 | 100.0 | 100.0 | 100.0 |
| 15 | 5 | 0.3 | 0.0 | 0.4 | 0.0 | 0.0 | 0.0 |
| 15 | 5 | 0.3 | 0.1 | 0.2 | 0.1 | 0.0 | 0.0 |
| 15 | 5 | 0.3 | 0.3 | 0.8 | 1.4 | 1.2 | 0.9 |
| 15 | 5 | 0.5 | 0.0 | 0.8 | 0.8 | 0.0 | 0.0 |
| 15 | 5 | 0.5 | 0.1 | 0.5 | 1.1 | 0.7 | 0.2 |
| 15 | 5 | 0.5 | 0.3 | 1.0 | 4.9 | 7.9 | 18.7 |
| 15 | 5 | 0.8 | 0.0 | 9.4 | 65.9 | 87.4 | 100.0 |
| 15 | 5 | 0.8 | 0.1 | 9.3 | 69.1 | 92.1 | 100.0 |
| 15 | 5 | 0.8 | 0.3 | 9.3 | 85.4 | 99.2 | 100.0 |
| 15 | 10 | 0.3 | 0.0 | 0.3 | 0.0 | 0.0 | 0.0 |
| 15 | 10 | 0.3 | 0.1 | 0.5 | 0.0 | 0.0 | 0.0 |
| 15 | 10 | 0.3 | 0.3 | 3.7 | 2.2 | 1.1 | 2.0 |
| 15 | 10 | 0.5 | 0.0 | 2.8 | 0.3 | 0.0 | 0.0 |
| 15 | 10 | 0.5 | 0.1 | 3.6 | 0.5 | 0.0 | 0.1 |
| 15 | 10 | 0.5 | 0.3 | 7.1 | 14.5 | 29.8 | 78.0 |
| 15 | 10 | 0.8 | 0.0 | 51.0 | 96.4 | 100.0 | 100.0 |
| 15 | 10 | 0.8 | 0.1 | 51.8 | 99.2 | 100.0 | 100.0 |
| 15 | 10 | 0.8 | 0.3 | 65.9 | 100.0 | 100.0 | 100.0 |

984 **Observed $D_s(\mathbf{R}_{\text{Pop}}, \mathbf{R}_{\text{Sm}})$ Values**

985 In addition to the estimated marginal means shown in the main text, the following
 986 figures (Figures B1–B4) show box-plots of $D_s(\mathbf{R}_{\text{Sm}}, \mathbf{R}_{\text{Pop}})$ for each condition in the
 987 simulation design. These box-plots match well with the estimated marginal means shown in
 988 the main text. However, notice that some conditions in these figures are missing box-plots
 989 (e.g., three factors, 15 subjects per item, 10 items per factor, $v_E = 0$, and Loading = 0.3)
 990 because no indefinite tetrachoric correlation matrices were produced for those conditions.

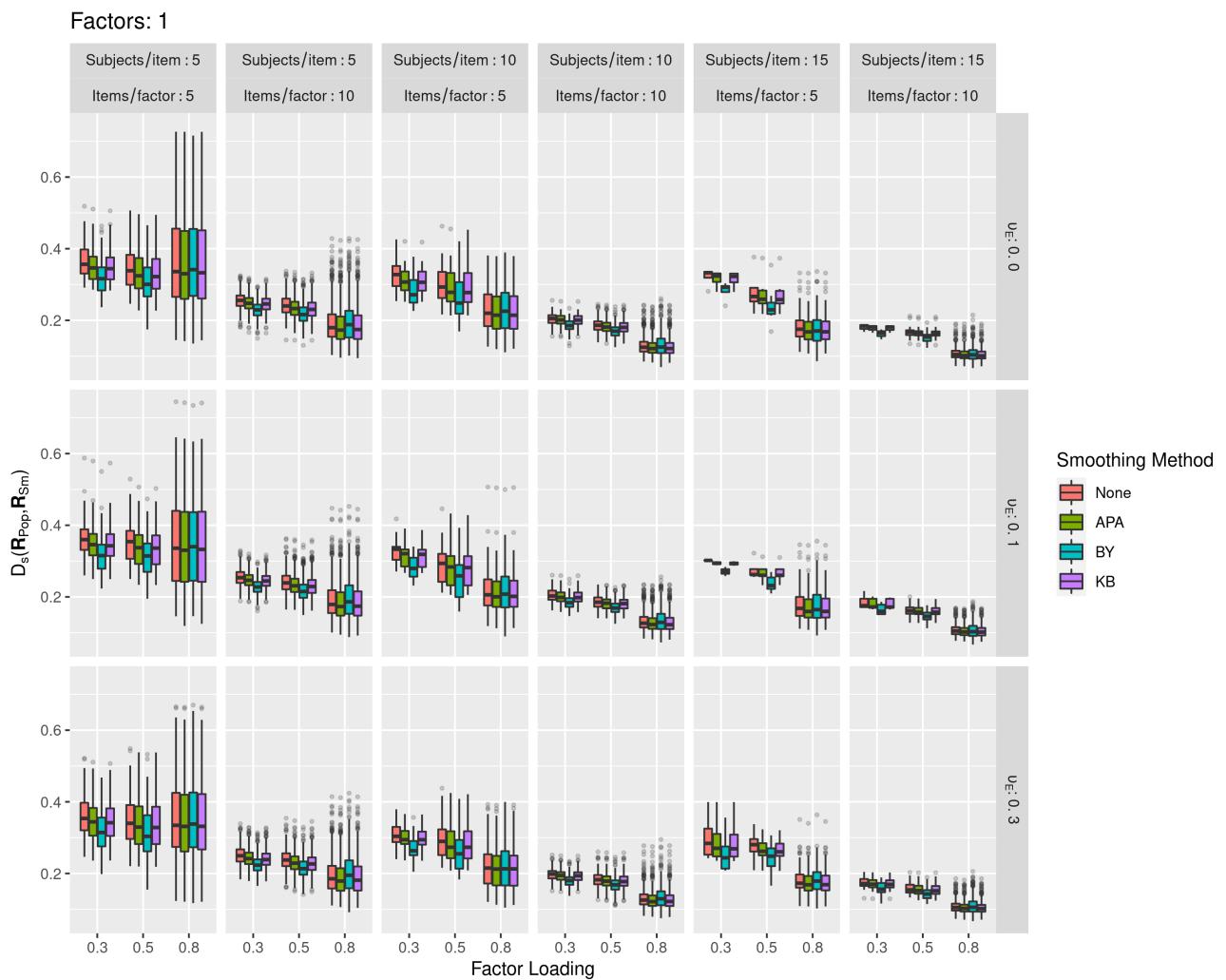


Figure B1. $D_s(\mathbf{R}_{\text{Pop}}, \mathbf{R}_{\text{Sm}})$ values for one-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

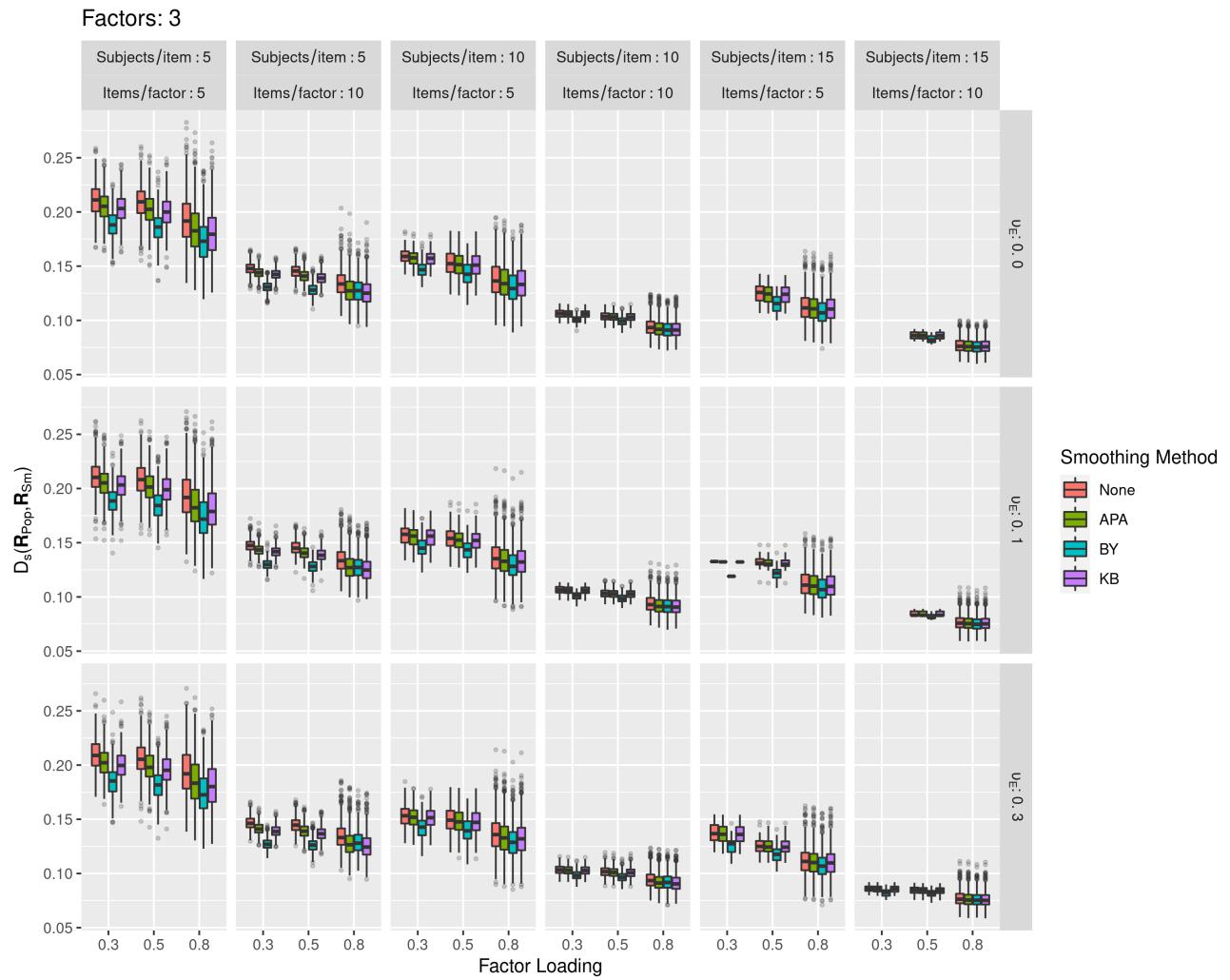


Figure B2. D_s(R_{Pop}, R_{Sm}) values for three-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

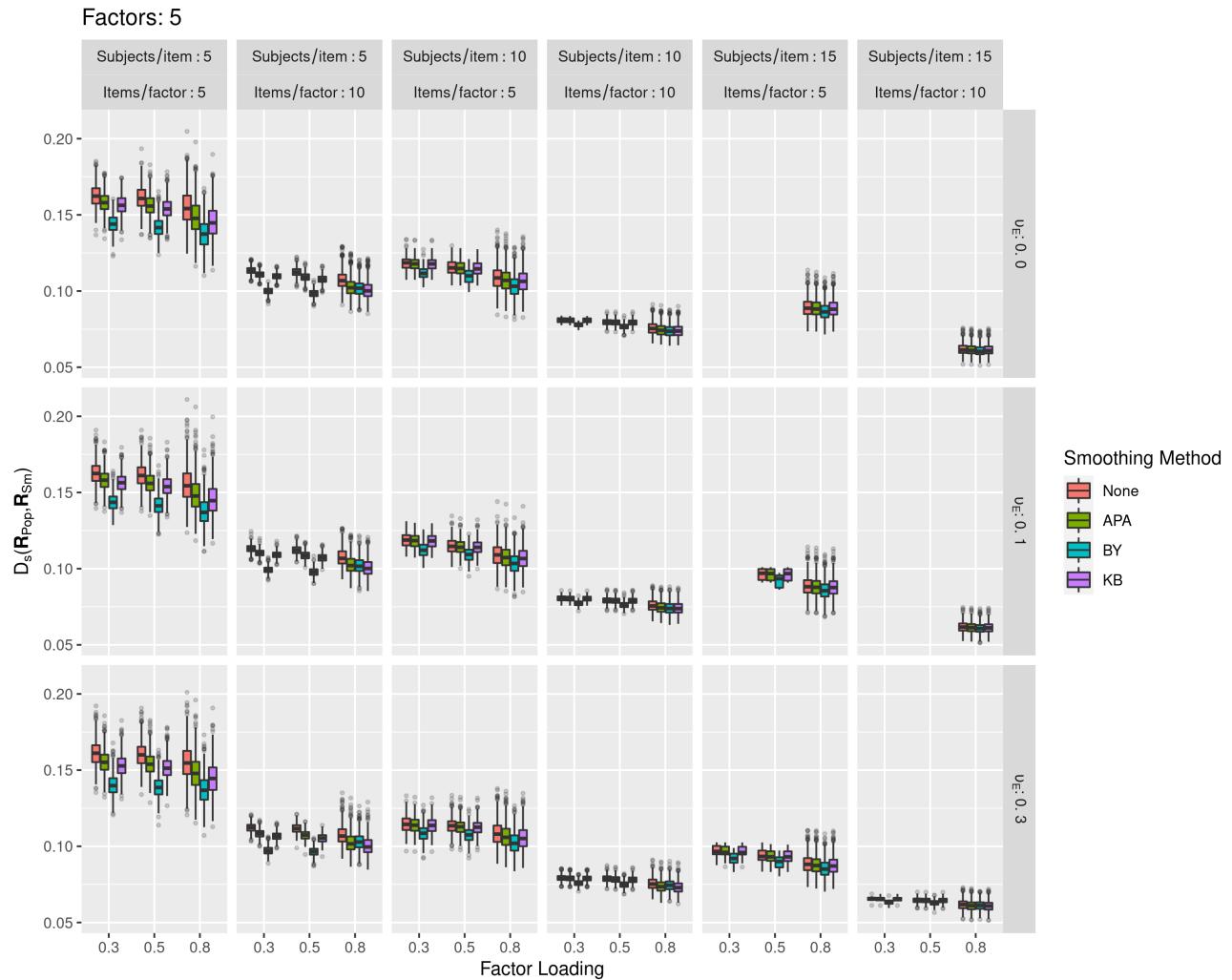


Figure B3. D_s(R_{Pop}, R_{Sm}) values for five-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

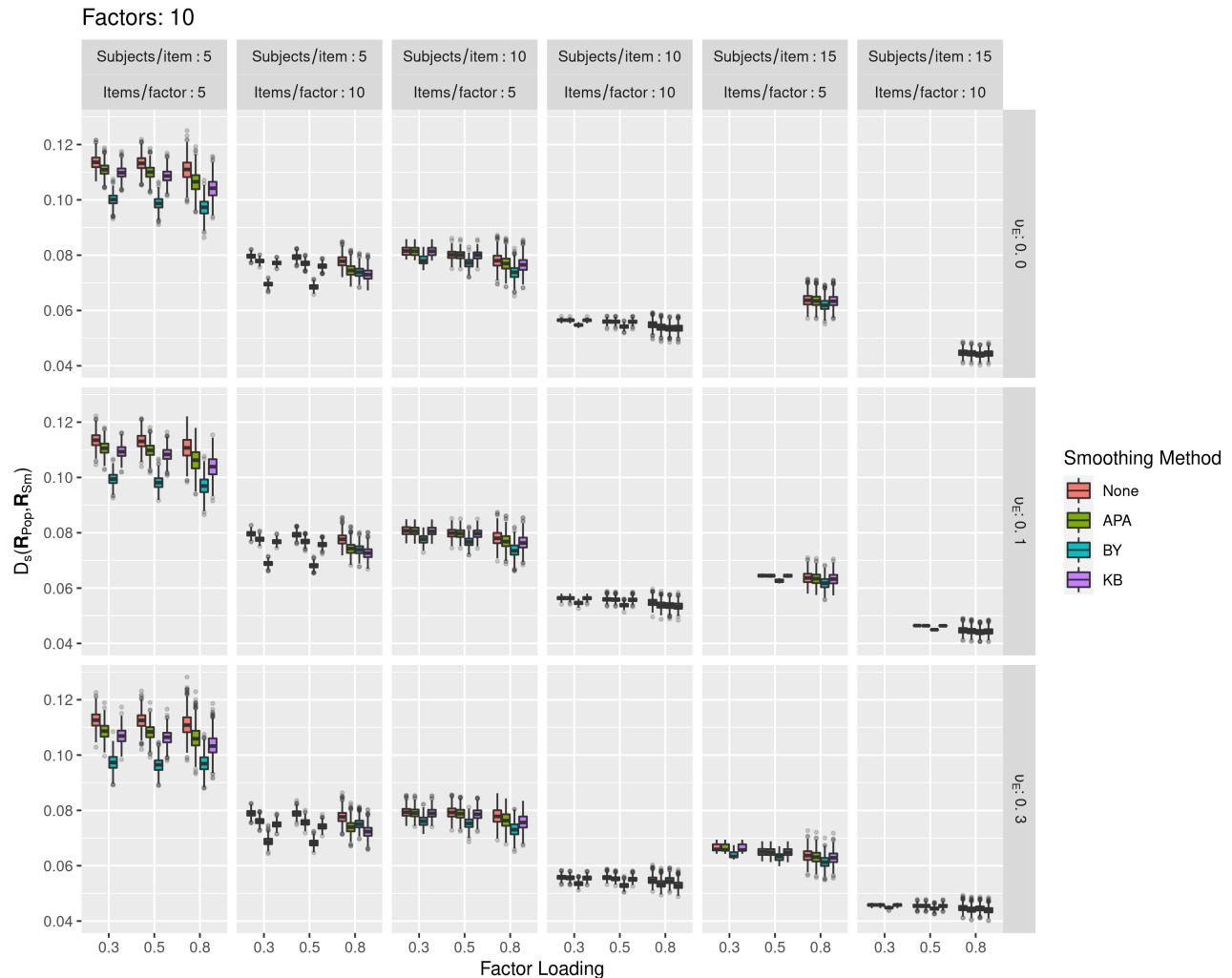


Figure B4. D_s(R_{Pop}, R_{Sm}) values for ten-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

991 Observed RMSE(\mathbf{F} , $\hat{\mathbf{F}}$) Values

992 Figures B5–B8 in this section show box-plots of $\text{RMSE}(\mathbf{F}, \hat{\mathbf{F}})$ for each condition in the
993 study design. Similar to the figures in the previous section, these box-plots for the most part
994 agree well with the estimated marginal means presented in the main text, but are missing
995 data for conditions with no indefinite tetrachoric correlation matrices.

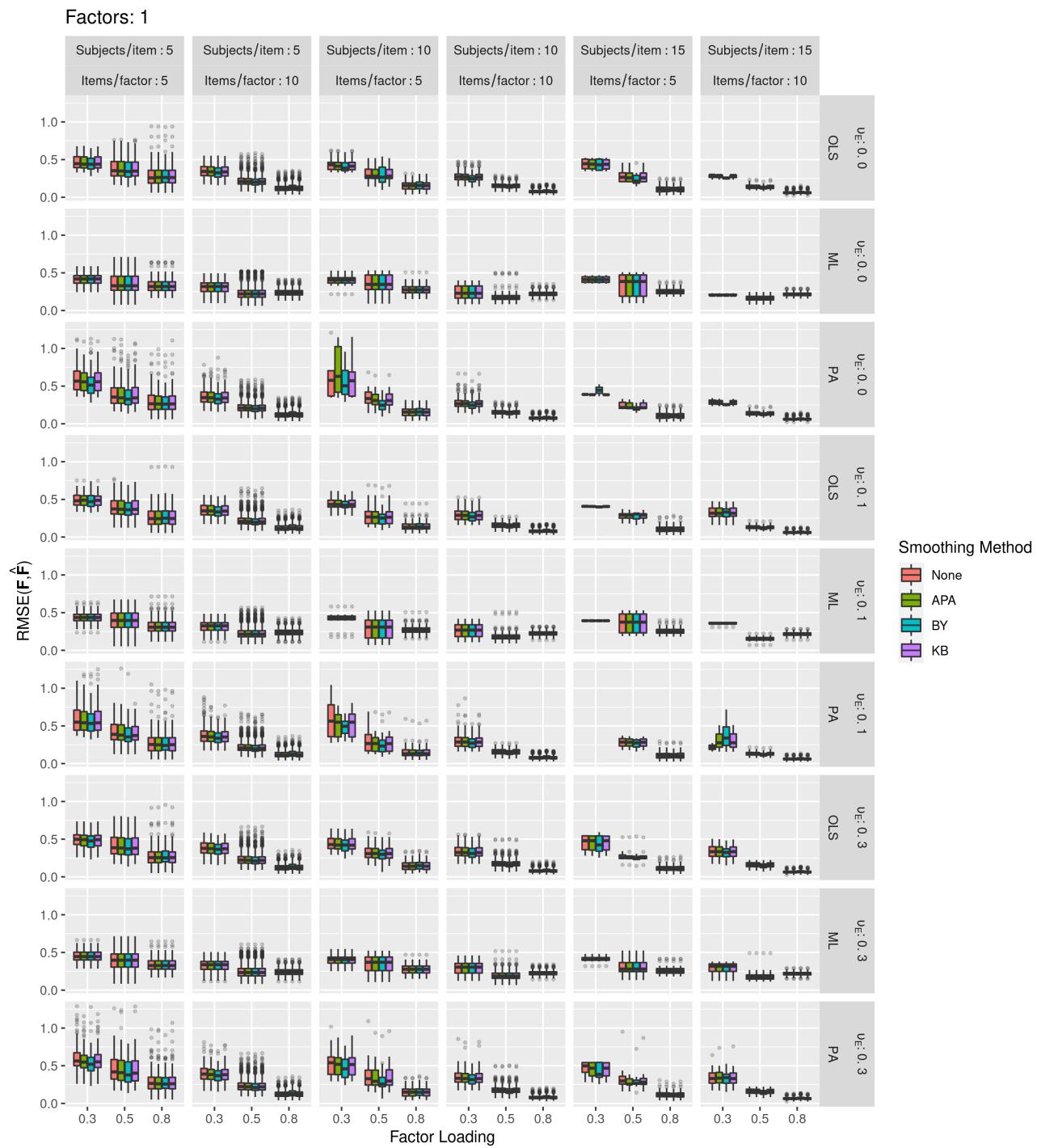


Figure B5. RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values for one-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; OLS = Ordinary least squares; ML = Maximum likelihood; PA = Principal axis; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

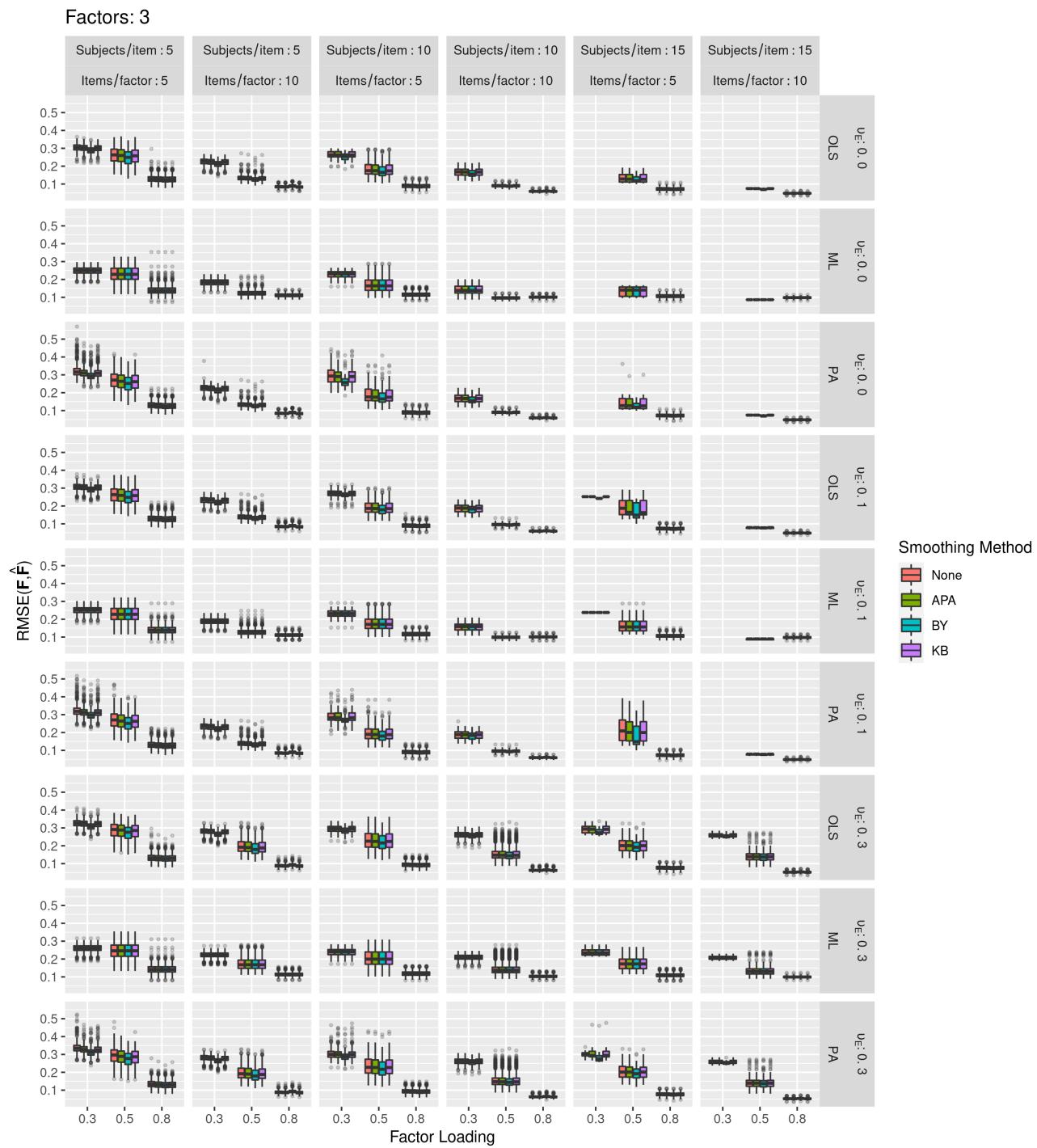


Figure B6. RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values for three-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; OLS = Ordinary least squares; ML = Maximum likelihood; PA = Principal axis; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

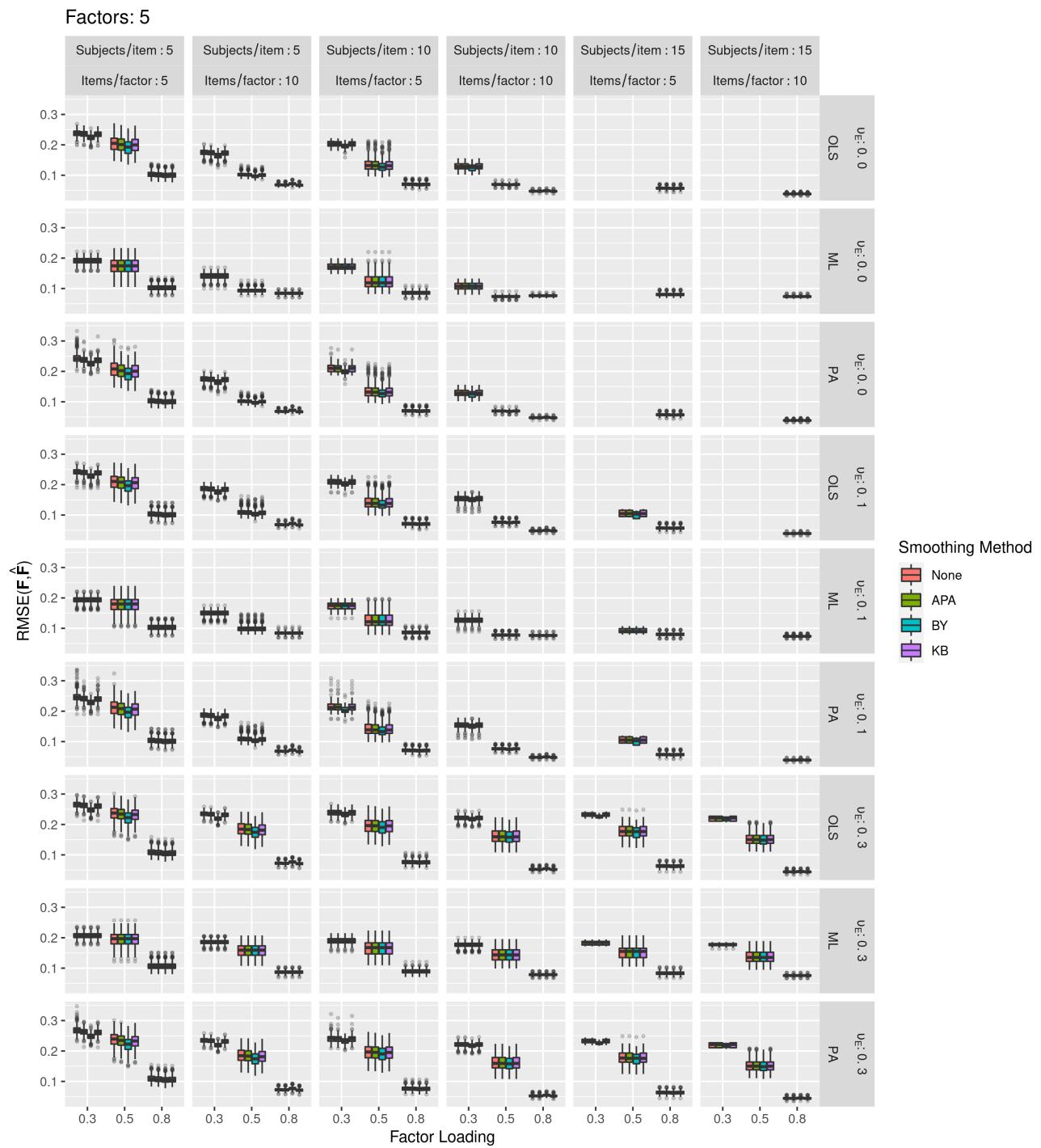


Figure B7. RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values for five-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; OLS = Ordinary least squares; ML = Maximum likelihood; PA = Principal axis; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.

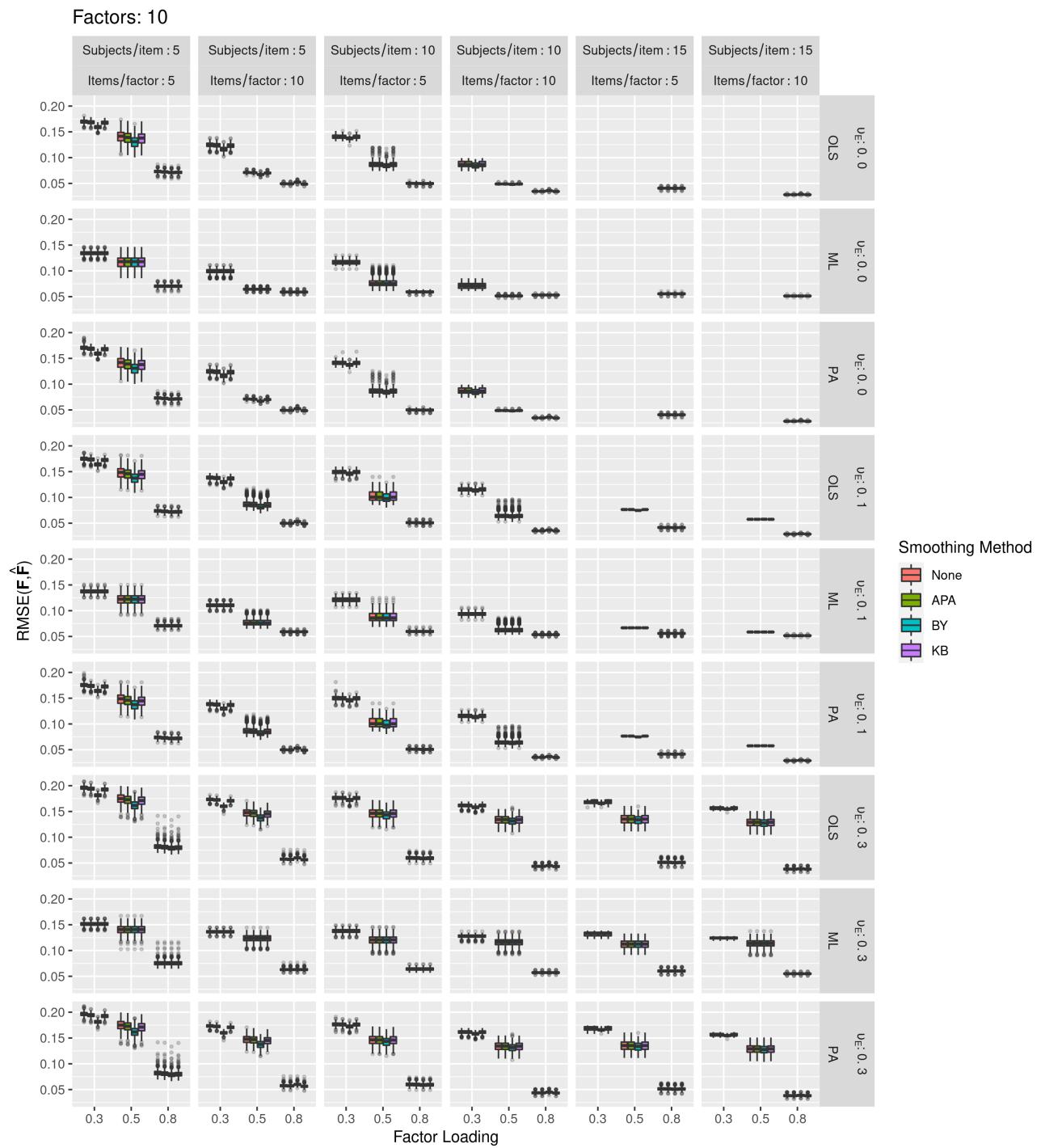


Figure B8. RMSE($\mathbf{F}, \hat{\mathbf{F}}$) values for ten-factor models. APA = Higham (2002); BY = Bentler-Yuan (2011); KB = Knol-Berger (1991); None = no smoothing; OLS = Ordinary least squares; ML = Maximum likelihood; PA = Principal axis; v_E = Proportion of uniqueness variance redistributed to the minor common factors representing model approximation error.