1

2

3

**Reducing the bias of norm scores in non-representative samples: Weighting as an**

**adjunct to continuous norming methods**

6

7

Sebastian Gary[1,] Alexandra Lenhard[1], Wolfgang Lenhard[2*], David S. Herzberg[3]

9

[1] Test Development Center, Psychometrica, Dettelbach, Bavaria, Germany

[2] Institute of Psychology, University of Würzburg, Bavaria, Germany

[3] WPS, Torrance, CA, USA

13

* Corresponding author:

E-Mail: wolfgang.lenhard@uni-wuerzburg.de

16

1

**Abstract**

17

18     In this simulation study, we investigated whether the accuracy of normed test scores

19  derived from non-demographically representative samples can be improved by combining

20  semi-parametric continuous norming methods with compensatory weighting at the raw score

21  level.  In a simulated reference population, we modeled a latent cognitive ability with a

22  typical developmental gradient, along with three demographic variables that were correlated

23  to varying degrees with the latent ability. We then simulated five additional populations

24  representing patterns of non-representativeness that might be encountered in real-world test

25  development research. We subsequently drew smaller normative samples from each

26  population, and used a one-parameter logistic IRT model to generate simulated test results for

27  each individual in the samples. Using these simulated data, we applied continuous,

28  regression-based norming techniques, both with and without compensatory weighting. The

29  weighting technique reduced the bias of the normed scores when the degree of non-

30  representativeness was relatively small. The technique was less effective with larger

31  departures from demographic representativeness, and when non-representativeness occurred

32  at extreme person locations. However, regardless of whether weighting was applied, the

33  observed norm errors resulting from demographic non-representativeness were small enough

34  to be ignored in most practical applications.

35

36  **Reducing the bias of norm scores in non-representative samples: Weighting as an**

37  **adjunct to continuous norming methods**

38      In the development of norm-referenced psychometric tests, demographically

39  representative samples provide the foundation for valid norm scores. An initial task for the

40  test developer is to identify those demographic variables that correlate most strongly with the

41  construct to be measured by the test. These variables typically include age, gender,

42  race/ethnicity, education level, and/or socioeconomic status. When measuring a developing

43  cognitive ability, age is the most important variable. Age has a stronger effect on test scores

44  than the other variables, especially when testing children and adolescents. Because of this,

45  Wechsler (1939, Chapter 3) recommended that same-age reference populations be used when

46  norming tests of intelligence and achievement.

47      Consequently, normative samples must be demographically representative not just

48  over the entire age range of the test, but also within smaller age groups containing individuals

49  who are at similar stages of development. Samples collected in this way create the basis for

50  *age-stratified* norms. Age stratification tends to increase the size of the normative sample

51  (along with its cost and the time needed to collect it). To mitigate the need for larger samples,

52  advanced mathematical methods have been developed to model the continuous relationship

53  between raw and normed scores across age (e.g., Gorsuch, 1983, quoted from Zachary &

54  Gorsuch, 1985; Cole, 1988; Cole & Green, 1992; A. Lenhard et al., 2018; W. Lenhard et al.,

55  2018; Stasinopoulos et al., 2018; A. Lenhard et al., 2019; W. Lenhard & Lenhard, 2021).

56      Some norm-referenced measures require stratification on variables other than age. For

57  example, measures of body mass index (BMI) need gender-stratified norms, because optimal

58  BMI for women is lower than that for men (Sang-Wook et al., 2015).

59      In other instances, it may be counter-productive to provide separate norms for the

60  different levels of a demographic variable. For example, some studies show that girls have

> **Commented [DH1]:** Can we find a more recent reference to support this point?

higher reading skills than boys (e.g., W. Lenhard et al., 2017; Price-Mohr & Price, 2017).

However, if a reading test is intended to identify those children who need additional support -

for example, children at the lowest decile of reading performance - then the use of gender-

stratified norms might result in biased outcomes. With gender-stratified norms, some girls

might be identified as needing additional support, even though they perform better on the test

than boys who are not identified as needing educational support.

**Addressing demographic imbalances: Stratification and post-stratification**

Besides creating stratified norms, several options exist for dealing with demographic

variables that are correlated with the latent ability being measured. An obvious course is to

increase the size of the normative sample. As sample size increases, the distribution of the

latent ability across gender (for example) increasingly approximates the distribution in the

reference population. However, cost and time constraints usually limit the size of the sample

available for norming.

A second approach is stratification, in which random sampling is conducted

independently within homogenous categories, or strata, defined by the demographic variable

(e.g., males, females). The goal is to have the category proportions in the normative sample

match, as closely as possible, the proportions in the reference population. For example, if

census data indicates that the reference population is composed of 50% males and 50%

females[1], then the researcher would sample males and females independently to match those

proportions in the normative sample.

However, it is not always possible to replicate population distributions through

stratified random sampling. One can randomly delete cases from over-represented strata, but

researchers are understandably reluctant to discard data. An alternative is to apply weighting

---

[1] For simplicity's sake, we assume at this point that the number of individuals who
identify as non-binary is negligible.

84    multipliers, or weights, to the data of individuals in the mis-represented strata. For example,

85    if a sample consists of 100 males and 50 females, a weighting multiplier of 2 could be applied

86    to the data obtained from females. Each test score from a female would then be treated as if

87    two females had obtained such a result. A weight $w_k$, assigned to an observation $x_i$ in

88    subsample $k$, thus indicates the number of individuals that this single observation represents.

89    The weights must therefore be calculated so that the proportion

90
$$p_k = \frac{w_k \cdot n_k}{\sum_k w_k \cdot n_k}$$

91    corresponds to the proportion of stratum $k$ in the reference population (with $n_k$ = size

92    of subsample $k$ in the norm sample). This weighting procedure is referred to as post-

93    stratification (Little, 1993; Park et al., 2004; Lumley, 2011, chapter 7).

94        Recently, Kennedy and Gelman (2021) recommended the use of multilevel regression

95    combined with post-stratification to correct for non-representative samples in studies of

96    psychological intervention. The authors suggest that weights can be used to adjust the means

97    of non-representative samples, to facilitate statistical comparisons among samples. However,

98    in constructing test norms from non-representative samples, the application of weights is

99    more complicated, because the norming process involves modeling both population means

100   and percentile ranks.

101       It is straightforward to take weights into account when calculating percentiles. As

102   described above, each test result is treated as if obtained by $w_k$ individuals. But, this simple

103   calculation runs the risk of introducing its own bias into the raw-to-norm-score relationship,

104   especially at the tails of the raw-score distributions. This risk occurs because the weights do

105   not change the variance of the distribution of raw scores within demographic subgroups, as

106   would be the case if more individuals were added to the subgroups. The potential distortion

107   of the variance of the raw score distributions increases with the magnitude of the weights

108   themselves. The risk for bias also *increases* as the number of individuals in a subgroup

**Commented [DH2]:** This equation not completely clear to me. I think the $k$s in the numerator need to be distinguished from the $k$s in the denominator.

5

109   *decreases* (as is expected at the tails of the raw-score distributions). Consequently, the

110   usefulness of weighting as a corrective procedure tends to diminish as the distributions of

111   demographic variables in the normative sample become increasingly divergent from those in

112   the reference population.

113        Because the potential distortions associated with weighting are most prominent at the

114   tails of the raw-score distributions, they can disproportionately affect the raw-to-norm-score

115   relationships for individuals of very high and/or very low ability. Unfortunately, these

116   extreme ability ranges are the ones where precise norm scores are most needed, because the

117   primary clinical applications of psychometric tests are to help diagnose disabilities, or,

118   alternatively, to identify gifted individuals.

119        As noted above, post-stratification is a method for dealing with normative samples

120   that are not representative, in term of demographic distributions, of the reference populations

121   from which they are drawn. An additional complicating factor is that the common

122   demographic variables of gender, socio-economic status (SES), race/ethnicity, and

123   geographic region are often inter-related, in terms of the effects they may have on test

124   performance. For example, areas with lower household income often have higher proportions

125   of non-white inhabitants. Because of such interactions, the most accurate approach to

126   stratification is to consider not only the marginal distributions of the demographic variables,

127   but also their cross-classifications, or joint distributions. In a complete crossing of the four

128   variables mentioned above, for instance, an individual could be classified as "female, low

129   SES, white, west region".

130        There are several practical difficulties with stratification based on the joint

131   distributions of demographic variables. For one, census data are often available only for

132   single demographic variables considered independently from each other, not for the cross-

133   classified categories of multiple variables. In addition, in one possible cross-classification of

134  gender, SES, race/ethnicity, and region, 192 joint cells (2 x 4 x 6 x 4) are created, some of

135  which require only a few individuals to meet census proportions. Collecting a sample that

136  meets these exacting specifications becomes a costly, lengthy process. In fact, with typical

137  sample sizes of 100 cases per age year in tests of cognitive ability (e.g., Kaufman &

138  Kaufman, 2004; Wechsler, 2008; Wechsler, 2014), it is not possible to replicate the census

139  proportions in every cross-classified cell, because some of the joint percentages specify less

140  than a single individual in a cell.

141  **Raking**

142      The raking procedure (Ireland & Kullback, 1968; Kalton & Flores-Cervantes, 2003) is

143  an approach to post-stratification that attempts to mitigate the practical challenges of

144  sampling based on a complete crossing of demographic variables. Raking does not draw on

145  the explicit joint distributions associated with all possible cross-classifications. Instead, the

146  post-stratification weights are determined in an iterative process based on the marginal

147  distributions of each demographic variable. That is, the weights assigned to the demographic

148  categories are adjusted successively and, if necessary, repeatedly, until they no longer

149  change. The procedure is termed "raking" because it is analogous to smoothing out the soil in

150  a garden bed by repeatedly raking in different directions. Studies have shown that the raking

151  procedure is convergent and delivers optimal asymptotically normal estimates for the joint

152  probabilities associated with a complete crossing of demographic variables (e.g., Ireland &

153  Kullback, 1968).

154      Although widely employed to correct for lack of representativeness in political polls

155  (Kalton-Flores-Cervantes, 2003), raking apparently has not been used in the norming of

156  psychometric tests, perhaps because it could introduce error into the raw-to-norm-score

157  relationships. As discussed above, demographic variables may interact with one another in

158  their effects on test scores, creating the need to consider the joint distributions of such

159  variables in developing norms. Because raking operates only on the marginal distributions

160  (i.e., it considers only the "main effects" of demographic variables on test scores), it may

161  magnify sources of error that stem from the interactions of these variables. However, these

162  potential risks remain at the level of speculation, because, to our knowledge, the effect of

163  raking on the accuracy of norm scores has never been studied.

164  **Effects of continuous norming on non-representativeness samples**

165      Continuous norming methods offer the advantage of using the properties of the entire

166  normative sample to correct local distributional anomalies in smaller subsamples (e.g., age

167  strata). Consequently, continuous norming methods may offer at least a partial remedy to

168  distortions caused by lack of demographic representativeness. The semi-parametric

169  continuous norming approach (SCN), first suggested by A. Lenhard and colleagues (A.

170  Lenhard et al., 2018; A. Lenhard et al., 2019; W. Lenhard & Lenhard, 2021), has been shown

171  to yield accurate norm scores with several non-optimal types of normative samples. One

172  advantage of SCN is that it does not make specific assumptions about distribution parameters,

173  and therefore can be applied to raw score distributions that are skewed, or that show floor

174  and/or ceiling effects. A second strength of SCN is that it performs better than parametric

175  approaches when applied to norm samples with the typical sample size of 100 per age cohort,

176  independent of the skewness of the raw score distributions (A. Lenhard et al., 2019).

177      In modeling the trajectories of percentile ranks across age groups, SCN (as

178  implemented in the cNORM package in R, A. Lenhard et al., 2018) *does not* rely on splines.

179  As a result, these trajectories are relatively rigid. As noted above, this feature of SCN

180  modeling tends to reduce the influence of error variance in local age groups, including that

181  caused by lack of demographic representativeness. As might be expected, therefore, SCN

182  generally produces less norming error than methods that determine raw-to-norm-score

183  mapping separately for each age group (W. Lenhard & Lenhard, 2021).

> **Commented [DH3]:** I don't understand what "rigid" means in this context. It seems you are arguing that because of this rigidity, error in certain age groups can be reduced by the influence of age-groups that have greater representativeness. What's missing is the logical link between rigidity of trajectories and the outcome of less error. Why does the former lead to the latter?

**Goals of the current simulation study**

184

185    We have described two techniques (SCN, raking) that may help ameliorate bias in

186  norm scores resulting from non-representative normative samples. To date, no research has

187  investigated whether the combination of both methods (in the following referred to as

188  *weighted continuous norming,* or WCN) further improves the accuracy of norm scores,

189  compared to SCN alone. A further open question is whether the mathematical

190  transformations wrought by SCN and raking might interact in a way that would *increase* the

191  bias of norm scores, within certain ability ranges.

192    The goal of the current study, therefore, was to evaluate the benefits and risks of

193  applying WCN[2] to non-representative normative samples. To this end, we conducted a

194  norming procedure on measure of a simulated cognitive ability that increases with increasing

195  age. Furthermore, we modeled the effects of three simulated demographic variables on the

196  cognitive measure. For convenience, we labeled the simulated demographic variables as

197  "education", "ethnicity", and "geographic region". We modeled education so that it would

198  have a stronger effect on the cognitive measure than ethnicity or region.

199    To provide input for the norming procedure, we generated six simulated population-

200  level data sets: a reference population that embodied the benchmark distributions of the three

201  demographic variables; and five non-representative populations, in which the distribution of

202  these demographic variables differed from the reference population. Each of these

203  populations has six equal-sized age cohorts. Table 1 summarizes the differences among the

204  six simulated populations.

205    Table 1: Simulated populations for norming input

---

[2] WCN incorporates SCN, as implemented in the R package cNORM (A. Lenhard et al., 2018).

| No. | Label | Description | Hypothesized effects on distribution of cognitive ability variable |
|---|---|---|---|
| 1 | Reference | Benchmark distributions of demographic variables; the standard of comparison for describing the "representativeness" of the other simulated populations. | Not applicable (benchmark population). |
| 2 | Mild under-representation of high education | Lower proportion of high-education individuals, higher proportion of low-education individuals, than Population 1. | Both mean and variance affected. |
| 3 | Moderate under-representation of high education | The pattern of divergence of education proportions is similar to population 2, but the degree of non-representativeness is greater. | Both mean and variance affected. |
| 4 | Under-representation of both low and high education | Both tails of the education distribution have lower proportions than Population 1. | Only variance affected. |
| 5 | Biased joint distributions | Marginal distributions of demographic variables match population 1; joint distributions (cross classifications) do not | Only variance affected. |

| | | | |
|---|---|---|---|
| | | match population 1. The pattern of non-representation alternates from over- to under-represented across the 27 (3 x 3 x 3) joint distributions. | |
| 6 | Clustered sampling | Marginal and joint distributions of demographic variables match population 1, but only when averaged across all six age cohorts. Within each age cohort, two-thirds of the joint distribution cells contain no data. | ??? (original manuscript does specify) |

**Commented [DH4]:** I would prefer the label "clustered distributions" for this population, because I want to keep a clean distinction between the generation of *population*-level data sets, and subsequent drawing of normative *samples*.

**Commented [DH5]:** What is the hypothesized effect on cognitive ability mean, variance for population 6?

206

207  Because raking incorporates only marginal distributions, we expect it to have little

208  effect in populations 5 (biased joint probabilities) and 6 (clustered sampling). In these two

209  populations, non-representativeness occurs only at the level of joint distributions (cross-

210  classifications), not at the level of marginal distributions.

211  For the current study, our pre-registered hypotheses were as follows:

**Commented [DH6]:** I would prefer to have the information about registration in a footnote, instead of in the main narrative.

212  1.We expected a main effect of norming method, such that WCN would lead to less-

213  biased estimates of the norm scores than SCN, where "bias" is quantified in terms of root

214  mean square error ($RMSE$) and mean signed difference ($MSD$).

215  2. We expected an interaction between norming method and the degree of non-

216  representativeness of the input data. Specifically, we expected that as the non-

217 representativeness of the normative sample increased, norm-score bias would increase for

218 both methods, but that the increase in bias would be smaller for WCN than for SCN.

219     3. We expected that the simple effect of WCN in reducing norm-score bias would

220 vary depending on person location on the cognitive variable. Specifically, we expected that

221 WCN would be less effective at reducing bias at the tails of the cognitive ability distribution

222 than in the central region of that distribution.

> **Commented [DH7]:** The last part of this section (which I deleted) is too vague to be part of a hypothesis; suggest simply reviewing findings post-hoc in results or discussion.

223 <div align="center">**Methods**</div>

224 *Overview*

225     Our study proceeded through the following steps:

226     1. Modeling a latent cognitive ability with a typical age-related growth curve.

227     2. Generating data sets for the reference population and five additional simulated

228        populations.

229     3. Drawing normative samples from each simulated population.

230     4. Generating simulated raw scores for a test of the cognitive ability.

231     5. Applying WCN and SCN to the raw scores from the normative samples.

232     6. Generating norm scores based on the reference population, as a standard of

233        comparison.

234     7. Testing the study hypotheses with ANOVAs, using RMSE and MSD as dependent

235        variables.

236 *Modeling cognitive ability*

237     To provide a basis for a modeled cognitive ability that develops with age, we

238 envisioned a reference population divided into six age cohorts, spanning one year each. We

239 conceptualized a cognitive ability that increases in each successive age group, as is typical

240 with cognitive development during childhood. We further specified that this cognitive ability

241 is influenced by the three demographic variables, each of which has three categories:

242        •   Education: low, medium, high

243        •   Ethnicity: native, mixed, non-native

244        •   Region: south, east, northwest

245 In broad terms, therefore, our model states that cognitive ability is a function of age and the

246 three demographic variables.

247        We operationalized the effect of the demographic variables on cognitive ability by

248 assigning three levels of mean cognitive ability (below average, average, above average) to

249 the three categories of each demographic variable, according to the matrix shown in Table 1.

250 **Table 1**

251 *Assignment of cognitive ability levels to demographic categories*

| Demographic variable | Below-average ability | Average ability | Above-Average ability |
|---|---|---|---|
| Education | low | medium | high |
| Ethnicity | native | mixed | non-native |
| Region | south | east | northwest |

252 Importantly, this mapping of ability level to demographic category remains constant in all

253 study conditions. Thus, by changing the distributions of demographic categories across

254 simulated populations, we simultaneously manipulate the distributions of cognitive ability.[3]

255        We then specified benchmark distributions for the demographic variables, which

256 would be enacted in the simulated reference population data set. The benchmark

257 demographic distributions must be understood in terms of a complete cross-classification of

258 the three demographic variables, which yields a 27-cell matrix with a 3 (low, medium, high

259 education) x 3 (native, mixed, non-native ethnicity) x 3 (south, east, northwest region)

---

[3] To limit the complexity of the simulation, we constrained the correlations among the demographic variables to zero.

260 structure. Table 2 follows this structure in specifying the benchmark demographic

261 distributions.

262 **Table 2**

263 *Distributions of demographic variables, by category, in the reference population*

264

|  |  | low education 40% | medium education 20% | high education 40% |
|---|---|---|---|---|
| ethnicity: native 30% | region: south 60% | 7.2% | 3.6% | 7.2% |
|  | region: east 20% | 2.4% | 1.2% | 2.4% |
|  | region: north-west 20% | 2.4% | 1.2% | 2.4% |
| ethnicity: mixed 40% | region: south 60% | 9.6% | 4.8% | 9.6% |
|  | region: east 20% | 3.2% | 1.6% | 3.2% |
|  | region: north-west 20% | 3.2% | 1.6% | 3.2% |
| ethnicity: Non-native | region: south 60% | 7.2% | 3.6% | 7.2% |
|  | region: east 20% | 2.4% | 1.2% | 2.4% |

| | | | | |
|---|---|---|---|---|
| 30% | region: north-west | | | |
| | 20% | 2.4% | 1.2% | 2.4% |

265

266    The benchmark marginal distributions of the demographic variables are shown in the

267    table margins (for education and ethnicity) and in the left-most column of the nested rows

268    (for region). The joint distributions for the complete cross-classification of the three variables

269    are shown in the table cells.

270    The structure of Table 2 provides a basis for understanding the demographic

271    manipulations that were applied to create five additional simulated populations. As described

272    previously, each cell of the table corresponds to a certain mean cognitive ability level, which

273    is determined by the demographic cross-classification of that cell. Thus, referring back to

274    Table 1, the cell in Table 2 with the highest mean cognitive ability is high education, non-

275    native ethnicity, northwest region, which appears in the lower-right corner of the table,

276    constituting 2.4% of the reference population. Conversely, the cell with the lowest mean

277    cognitive ability is low education, native ethnicity, south region, which appears in the upper-

278    left corner of Table 2, constituting 7.2% of the reference population.

279    Within the reference population, each row of data includes age, level of cognitive

280    ability, and classifications on education, ethnicity, and region. The values for the demographic

281    classifications are assigned according to the percentages in Table 2. The row-wise values of

282    cognitive ability are based on a distinct mean value for each cell[4] of Table 2. This cell-wise

283    mean is calculated by the following polynomial equation:

---

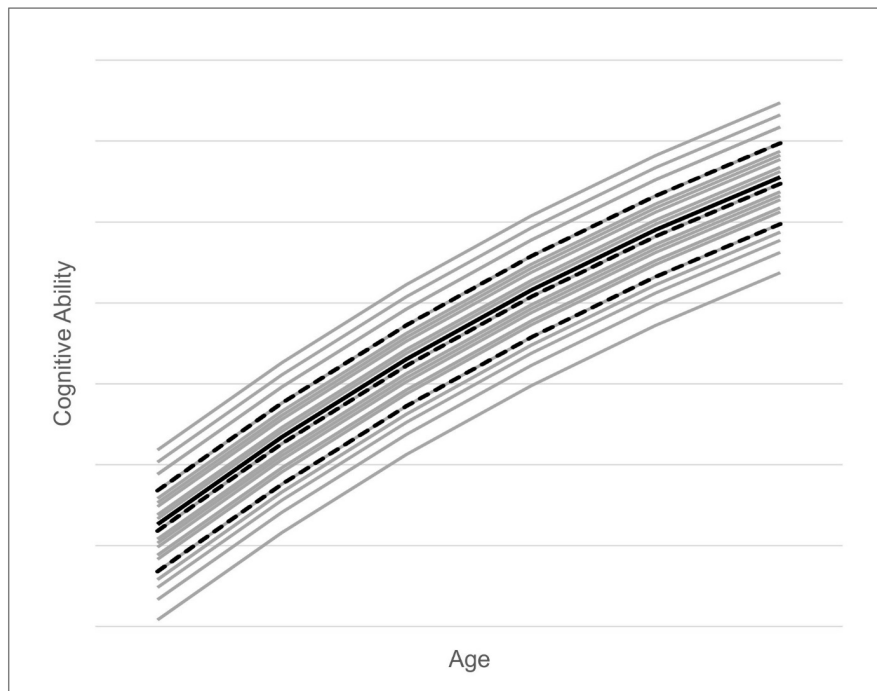[4] *SD* is constrained to 1 across all cells of Table 2.

15

284 $$M_{(age, education, ethnicity, region)} = -1.5 \cdot education - 0.25 \cdot ethnicity -$$

285 $$0.1 \cdot region - 0.05 \cdot ethnicity \cdot region + 1.2 \cdot age - 0.06 \cdot age^2 + 0.0001 \cdot age^4$$

286 (1)

287 As a result of this equation, each demographic variable exerts a different effect on

288 cognitive ability:

289 • Education: correlates at r = -.78 with cognitive ability (large effect)

290 • Ethnicity: $r$ = -.54 (medium effect)

291 • Region: $r$ = -.31 (small effect)

292 Figure 1 provides a graphic depiction of the modeled cognitive ability in the reference

293 population.

294

295 *Figure 1.* Modeled cognitive ability in reference population

296

297

298

299

300     The figure shows mean cognitive ability increasing across the six age cohorts. The

301     solid black line represents the reference population mean. The dashed black lines represent

302     the marginal mean cognitive abilities for the low, medium and high categories of education,

303     the demographic variable with the largest effect on cognitive ability. The grey lines represent

304     the mean cognitive abilities associated with the 27 demographic cross-classifications. The

305     highest grey line is high education, non-native ethnicity, northwest region, the cell of Table 2

306     with the highest mean cognitive ability. The lowest grey line is low education, native

307     ethnicity, south region, the cell of Table 2 with the lowest mean cognitive ability.

308     *Generation of simulated population data sets and normative samples*

309     To generate the reference population data set, we drew 24 million[5] pairs of random

310     numbers (4 million per age cohort), each pair representing one individual. The first random

311     number was uniformly distributed between 0 and 6 and represented age in years. The second

312     number was normally distributed with $M = 0$ and $SD = 1$ and represented the cognitive ability

313     of the individual with respect to other individuals of the same stratum and age. This random

314     number was converted into the specific cognitive ability value for an individual by adding the

315     mean cognitive ability for that individual's demographic cross-classification status (see Table

316     2 and Formula 1). Additionally, we $z$-standardized each cognitive ability value ($x_1$), using the

317     reference population mean $\mu$ and standard deviation $\sigma$ in formula 2,

---

[5] The reference population size is roughly based on the number of persons in the U.S.
population whose ages are within a span of six consecutive years (e.g., ages 0 to 5).

17

**Commented [DH8]:** Suggest labelling the age cohorts (0-5?) on the figure, for clarity.

**Commented [DH9]:** I think we need to say explicitly that the size of the population reference sample was aligned to the US population for ages 0-6 (if that is in fact what was meant). Or, "0-6" not supposed to represent actual ages, but merely serve as a means for classifying individuals into cohorts (in other words, it's nominal)?

**Commented [DH10]:** Do you really mean with respect to the same stratum AND age? Or just with respect to the same age cohort? The latter makes more sense to me, because the stratum is taken into account in the next step, when the random number is added to the stratum specific mean from Formula 1.

**Commented [DH11]:** This needs to be clarified. Do you mean that the second random number was nested within the first, such that there were six separate random draws for the second number, one for each age cohort?

$$\theta_{pop} = \frac{x_1 - \mu}{\sigma} \qquad (2)$$

where $\theta_{pop}$ represents an individual's location on the cognitive ability variable with respect to the entire reference population.

As noted previously, each individual in the reference population was assigned values on the demographic variables, such that marginal and joint distributions of these variables would match the distributions shown in Table 2. We then generated five additional simulated population data sets, using the same method described at the outset of this section. These additional simulated populations represented various violations of demographic representativeness that might be encountered in collecting normative data for the development of a psychometric tests The distributions of the demographic variables in these five additional data sets differed from the reference population as follows:

- Simulated Population 2: Mild under-representation of high education. The high education category was underrepresented (28% instead of 40%) and the low education category was overrepresented (52% instead of 40%). This manipulation affected both the mean and the variance of the cognitive ability variable.

- Simulated Population 3: Moderate under-representation of high education. The pattern of misrepresentation was the same as Population 2, but the degree of misrepresentation was greater (high education was 20% instead of 40%; low education was 60% instead of 40%). This manipulation affected both the mean and the variance of the cognitive ability variable.

- Simulated Population 4: Under-representation of both low and high education. Medium education was overrepresented (40% instead of 20%), and high and low education were underrepresented (30 % instead of 40 %). This manipulation attenuated the variance of the cognitive ability variable, but its mean was not affected.

342    • Simulated Population 5: Biased joint distributions. The joint distributions of the

343      demographic variables were varied from the reference percentages shown in Table 2,

344      such that:

345          o  Some of the joint cells were overrepresented, while some were

346             underrepresented.

347          o  The marginal distributions were identical to those in the reference population.

348      This manipulation increased the variance of the cognitive ability variable, but its

349      mean was not affected.

350    • Simulated Population 6: Clustered sampling. Within each age cohort, two-thirds of

351      the 27 demographic cross-classification cells contained no data. In the remaining one-

352      third of cells, the number of individuals was tripled. This manipulation was applied to

353      different subsets of cells across age cohorts, such that when cell proportions were

354      summed across all age cohorts, the marginal and joint distributions of the

355      demographic variables were identical to the reference population.

356      In the five additional simulated population data sets, the cognitive variable was $z$-

357      standardized using Formula 2. Importantly, the values of $\mu$ and $\sigma$ were those from the

358      reference data set (Population 1), not from the data set whose values were being standardized.

359      From each of the six simulated populations, we drew 100 random samples of 600

360      individuals (100 cases per age cohort). These samples served as input to the norming

361      procedures.

362    ***Simulation of test results***

363      Using the one-parameter logistic (1-PL) model, we simulated a 31-item test to

364      generate test results for each individual in the normative samples. The 31 item difficulties ($\delta$)

365      were drawn randomly from a uniform distribution ranging from -3 and +3. The set of item

366      difficulties covered a range of about 3.7 standard deviations ($M = -0.04$, $SD = 1.64$), therefore

**Commented [DH12]:** Do we need to specify the details of this manipulation in a footnote (probably too long to include in the main narrative)?

**Commented [DH13]:** See my previous comment - I prefer the label "Clustered Distriutions".

**Commented [DH14]:** For this population, how are the mean and variance of the cognitive variable affected? Also, as with Population 5, do we want to provide more details about the manipulation in a footnote?

**Commented [DH15]:** We need a sentence explaining why it was done this way, that is, why use the reference mean and variance for $z$-standardization of the cognitive variables in the five additional simulated populations?

367    spanning a wide range of latent ability. The probability $p_{k,i}$ that an individual $k$ with the $z$-

368    standardized latent ability $\theta_{pop\_k}$ succeeded on item $i$, with difficulty $\delta_i$, was given by the

369    following 1-PL equation:

370    $$p_{k,i}(x_i = 1|\theta_{pop\_k}, \delta_i) = \frac{\exp(\theta_{pop\_k} - \delta_i)}{1+\exp(\theta_{pop\_k} - \delta_i)} \tag{3}$$

371    For every individual $k$ and item $i$ a uniformly distributed random number between 0 and 1

372    was drawn and compared to $p_{k,i}$. If $p_{k,i}$ exceeded the random number, the item was scored 1,

373    otherwise it was scored 0. Finally, each individual's scores on all 31 items were summed to

374    yield a raw total score on the simulated test.

375

376    *Application of weighted and unweighted norming procedures*

377    For each raw score in the normative samples, we applied WCN and SCN to generate

378    IQ-type standard scores (M = 100, SD = 15) for each norming method. These scores were

379    labeled $IQ_{WCN}$ and $IQ_{SCN}$. For both WCN and SCN, these IQ scores were calculated with

380    *cNORM*, an R package that employs continuous norming (A. Lenhard et al., 2018). Weights

381    were not used for SCN.

382    To calculate the weights for WCN, we used the *rake* function from *survey* (Lumley,

383    2011), an R package that implements the raking procedure described earlier in this

384    manuscript. Additionally, we standardized the weights to make them easier to interpret. We

385    divided each weight by the smallest weight in the respective norm sample, thereby setting the

386    weight of the most overrepresented group in the sample to 1.[6]

387    Using weights required modifications to the standard cNORM functions. In WCN,

388    weights are applied initially in the ranking procedure, where each raw score is assigned a

---

[6] This transformation does not affect the distribution of the weights.

---

**Commented [DH16]:** The original manuscript had "underrepresented" here, but I believe you meant overrepresented, because it is the latter groups that have the smallest weights.

389   percentile rank. Weighting requires that raw scores from certain demographic groups be

390   counted more than once in this ranking process. For example, if a certain group is assigned a

391   weight of 2, each raw score from that group must be counted twice during ranking. Because

392   the weights, as initially calculated by the *rake* function, are floating point numbers (not

393   integers), they were multiplied by $10^6$ and then rounded to whole numbers prior to ranking.

394   This transformation made it possible, during ranking, to count each test score the number of

395   times indicated by its weight. Because of the high number of ties, the average rank was used

396   for further processing, following the usual cNORM procedures.[7]

397        In WCN, weights are also entered in cNORM's regression-based modeling procedure.

398   To perform the regression, cNORM draws on the *regsubsets* function of *leaps*, an R package

399   (Lumley, 2017). *regsubsets* includes the capacity to process weights in the regression

400   analysis.

401   ***Generating norm scores from the reference population***

402        To test the study hypotheses, we created a measure ($IQ_{best}$), in the same metric as

403   $IQ_{WCN}$ and $IQ_{SCN}$, which represented the "actual" person location on the cognitive ability

404   variable.  $IQ_{best}$ was derived from the distribution of raw scores in the entire reference

405   population (in contrast to $IQ_{WCN}$ and $IQ_{SCN}$, which were derived from the smaller normative

406   samples).

407        To compute $IQ_{best}$, we generated raw scores on the 31-item simulated test for the 24

408   million individuals in the reference population, using the previously described method. We

409   then partitioned the reference population by age, creating 365 equal-sized groups within each

410   of the six age cohorts. Each of the resulting 2190 age-groups consisted of about 11,000

411   individuals with the same "birthday". The raw scores were ranked and converted into IQ

---

[7] For a detailed description of the *cNORM* norming process, see A. Lenhard, Lenhard, & Gary, 2019.

412  scores using rank-based inverse normal transformation within each age group. As a result,

413  each row in the reference population data set included values for age, raw score and $IQ_{best}$.

**Hypothesis testing with RMSE and MSD**

415  As noted above, we drew 100 normative samples ($N = 600$) from each of the six

416  simulated population data sets. We conducted ANOVAs to test the study hypotheses in each

417  of these 600 normative samples. The ANOVAs compared $IQ_{best}$ to $IQ_{WCN}$ and $IQ_{SCN}$,

418  respectively, with *RMSE* and *MSD* as dependent variables. Both RMSE and MSD are

419  quantified in terms of IQ points.

420  *RMSE* is a summary measure of norming model error that includes both fixed and

421  variable error components (Lenhard & Lenhard, 2021). It was computed using the following

422  formula:

$$RMSE = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(IQ_. - IQ_{best})^2}, \qquad (4)$$

424  where $n$ is the number of cases and $IQ_.$ stands for either $IQ_{WCN}$ or $IQ_{SCN}$.

425  *MSD* is a measure of the tendency for a norming model to overestimate ($MSD > 0$) or

426  underestimate ($MSD < 0$) the actual person location. The formula used to calculate the *MSD*

427  was:

$$MSD = \frac{1}{n}\sum_{i=1}^{n}(IQ_. - IQ_{best}). \qquad (5)$$

429  To test Hypothesis 3, we divided the distributions of $IQ_{best}$, $IQ_{WCN}$ and $IQ_{SCN}$ into 11

430  intervals of 7.5 IQ points each. *RMSE* and *MSD* were calculated separately for each of these

431  intervals.

432  In general, the analytic approach was to conduct 6 (simulated population) x 11 (IQ

433  range) x 2 (norming method) mixed ANOVAs. Population was a between-groups factor, and

434  IQ range and norming method were within-groups factors. Because of size of the simulated

435    data sets, statistical power was high, and therefore the level of significance was set to $p = .01$.

436    The assumption of sphericity was tested and, where indicated, degrees of freedom were

437    corrected. Additionally, partial $\eta^2$'s were computed as measures of effect size. We further

438    specified that, in the norm score comparisons of interest, differences of less than 0.5 IQ were

439    too small to have any practical relevance.

> **Commented [DH17]:** How was this threshold determined?

440    **Results**

441    As indicated by Mauchly's test, sphericity assumptions were generally violated both

442    for *RMSE* and *MSD*. Therefore, degrees of freedom in all ANOVAs were corrected according

443    to the Greenhouse-Geisser method.

444    The 6 x 11 x 2 mixed ANOVAs yielded significant results for all main effects and

445    interactions ($p < .001$). We focus here on the effects that are most salient for testing our

446    hypotheses.

447    **Hypothesis 1: Main Effect of Norming Method**

448    The first hypothesis proposed that WCN would yield lower levels of norm-score bias

449    than SCN. This hypothesis was supported by tests of the main effects of norming method,

450    *RMSE*: $F(1, 594) = 94.93$, $p < .001$, $\eta^2 = .14$, *MSD*: $F(1, 594) = 3397.28$, $p < .001$, $\eta^2 = .85$.

451    *RMSE* was smaller for WCN ($M = 2.18$, $SE = .02$) than for SCN ($M = 2.36$, $SE = .02$). The

452    same was true for MSD (WCN: $M = 0.74$, $SE = .03$; SCN: $M = -0.24$, $SE = .03$).

> **Commented [DH18]:** These numbers appear to contradict the claim that MSD was smaller for WCN than SCN.
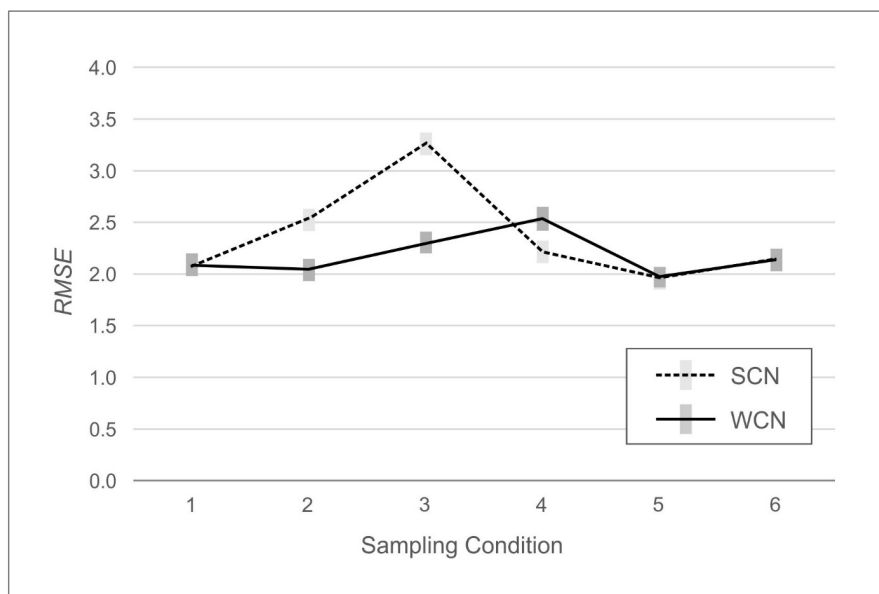
453    However, the analysis also detected significant interactions between norming method

454    and simulated population, indicating that the effects of weighting on norm-score bias varied

455    among the simulated populations, *RMSE*: $F(5, 594) = 98.98$, $p < .001$, $\eta^2 = .45$, *MSD*: $F(5,$

456    $594) = 764.77$, $p < .001$, $\eta^2 = .87$. As can be seen in Figure 2 (*RMSE*) and Figure 3 (*MSD*), in

457    two of the six simulated populations, weighting reduced bias in the normed scores. In

458    populations 2 (mild under-representation of high education) and 3 (moderate under-

459    representation of high education), *RMSE* was lower with WCN than with SCN, by an average

> **Commented [DH19]:** You'll have to edit the figures to reflect the change in nomenclature from "sampling condition" to "simulated population". Please check the other figures to see if they need similar changes.
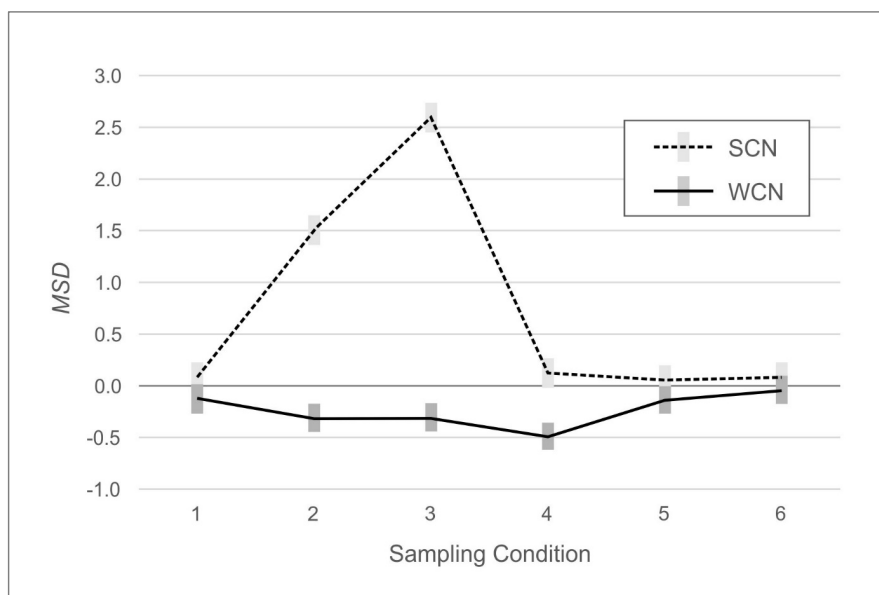
> **Commented [DH20]:** Please add to the figure caption, or in the main narrative, that these plots show the average RMSE and MSD, average over the 600 norm draws per simulated population.

460     of 0.48 IQ points and 0.97 IQ points, respectively. In populations 1 (reference), 5 (biased

461     joint probabilities) and 6 (clustered sampling), the difference in *RMSE* between WCN and

462     SCN approximated zero. In population 4 (under-representation of both low and high

463     education), WCN returned higher average *RMSE* than SCN, but the difference of

Figure 2. *RMSE* across simulated populations, with (WCN) or without (SCN) weighting. The grey rectangles represent 95% confidence intervals.



Figure 3. *MSD* across simulated populations, with (WCN) or without (SCN) weighting. The grey rectangles represent 95% confidence intervals.

468     0.32 IQ points was below the threshold of practical relevance.

469        The analysis of the *MSD* yielded similar results. In populations 2 (mild under-

470     representation of high education) and 3 (moderate under-representation of high education),

471     *MSD* was closer to the ideal value of zero for WCN (populations 2 and 3: -0.32 IQ points)

472     than for SCN (population 2: 1.51 IQ points; population 3: 2.59 IQ points). In populations 1

473     (reference), 5 (biased joint probabilities) and 6 (clustered sampling), *MSD* approximated zero,

474     regardless of the norming method. In population 4 (under-representation of both low and high

475     education), *MSD* deviated more from zero for WCN (-0.49 IQ points) than for SCN (0.13 IQ

476     points). As with *RSME*, these latter differences did not meet the criterion for practical

477     significance.

478     **Hypothesis 2: Interaction Between Norming Method and Degree of Non-**

479     **Representativeness**

480        Hypothesis 2 specified that as the non-representativeness of the normative samples

481     increased, norm-score bias would increase for both methods, but that the increase in bias

482     would be smaller for WCN than for SCN. To address this hypothesis, we compared

483     populations 2 and 3. Both populations were characterized by under-representation of the high

484     education group, but the magnitude of under-representation was greater in population 3 than

485     in population 2. Therefore, we performed two additional ANOVAs that limited the levels of

486     the between-groups factor to populations 2 and 3. Both analyses yielded a significant

487     interaction between norming method and simulated population, *RMSE*: $F(1, 198) = 26.01$,

488     $p < .001$, $\eta^2 = .12$, *MSD*: $F(1, 198) = 242.36$, $p < .001$, $\eta^2 = .55$.

489        The results of these analyses are visualized in Figures 2 and 3. The plots show the

490     interaction: *RMSE* and *MSD* are greater in magnitude in population 3 (moderate under-

491     representation) than in population 2 (mild under-representation), for both norming methods,

492     but the magnitude of increase in the error metrics is greater for SCN than for WCN.

26

493      Considering WCN in isolation, average *MSD* was approximately equal for both populations

494      (-0.32 IQ points). Average *RMSE*, on the other hand, did increase significantly in population

495      3, $F(1, 198) = 10.12$, $p = .002$, $\eta^2 = .05$.[8]

496      **Hypothesis 3: Effectiveness of WCN depends on person location**
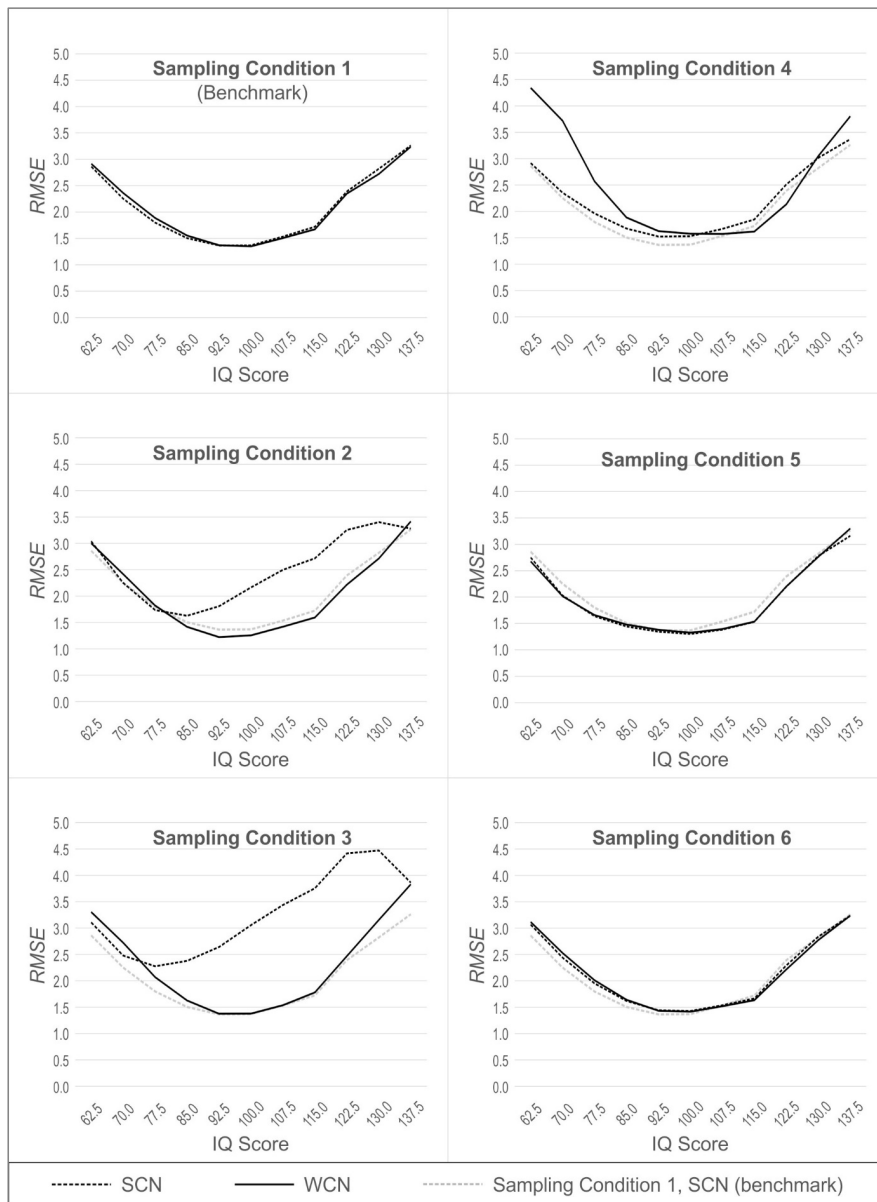
497          Hypothesis 3 proposed that WCN would be less effective at reducing bias at the tails

498      of the cognitive ability distribution than in the central region of that distribution. We tested

499      this hypothesis with two analytic approaches. First, we conducted 11 x 2 ANOVAs with

500      person location and norming method (WCN vs. SCN) as within factors, and *RMSE* and *MSD*

501      as dependent variables. We then examined how the effects of person location varied among

502      the simulated populations. We compared the performance of WCN in populations 2, 3, 4, 5

503      and 6 (which yield demographically non-representative normative samples, as described

504      earlier) to SCN in population 1 (which yields demographically representative normative

505      samples). SCN in population 1 therefore represents a benchmark condition, against which the

506      performance of WCN in the other non-representative populations can be measured. These

507      latter analyses used 11 x 6 ANOVAs, with simulated population as a between-groups factor.

508          The results of these analyses are illustrated in Figures 4 (*RMSE*) and 5 (*MSD*).

509      Because of the large number of comparisons, we report effects only if at least one of the

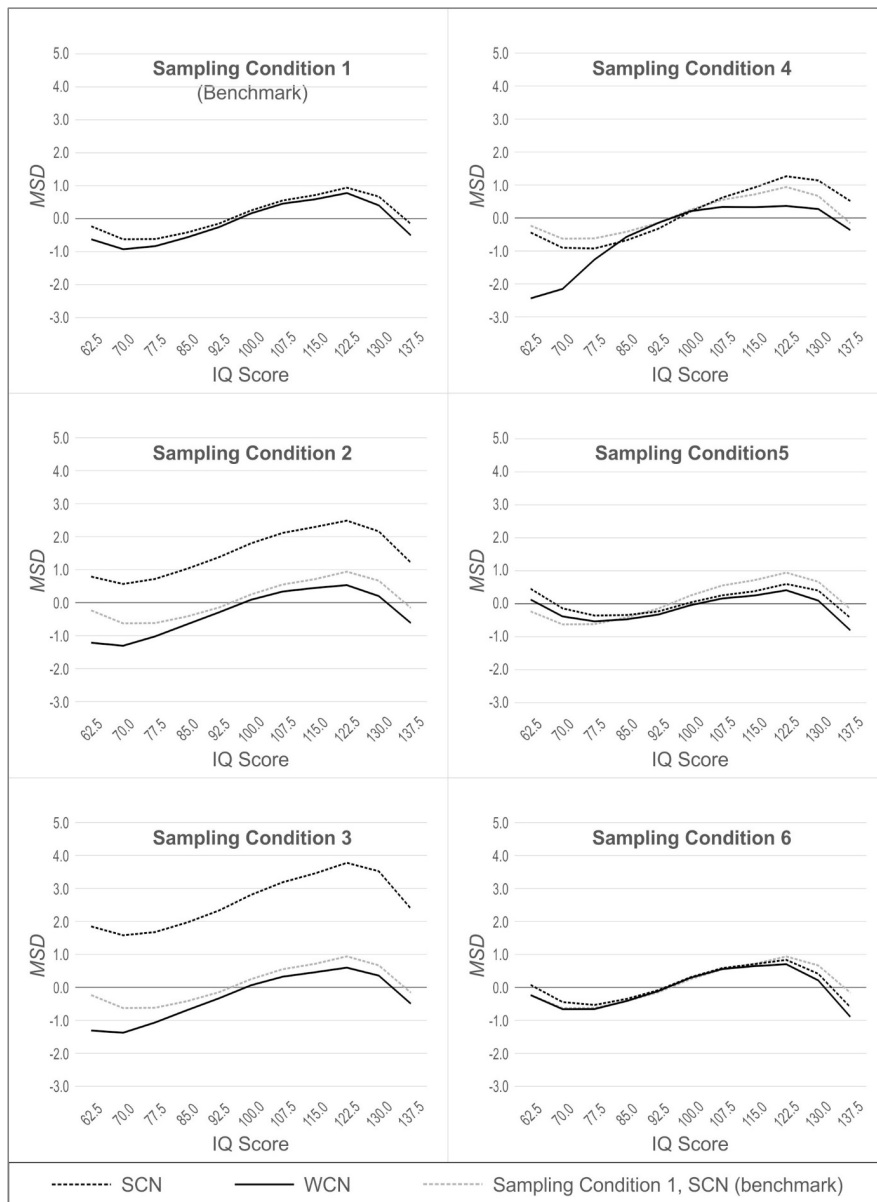510      differences within an analysis exceeded 0.5 IQ points.

511

---

[8]Average *RMSE* was 2.04 IQ points in sampling condition 2 and 2.29 IQ points in sampling condition 3.

**Figure 4.** *RMSE* across simulated populations, with (WCN) or without (SCN) weighting, as a function of person location. The dotted grey line represents SCN with norm samples drawn from Population 1 (benchmark).

515 **_Figure 5._** _MSD_ across simulated populations, with (WCN) or without (SCN) weighting, as a function of person
516 location. The dotted grey line represents SCN with norm samples drawn from Population 1 (benchmark).

517  *Population 1: Reference*

518       In normative samples drawn from the reference population, both ANOVAs yielded a

519  significant main effect of person location, *RMSE*: $F(2.68, 264.87) = 70.68$, $p < .001$, $\eta^2 = .42$,

520  *MSD*: $F(2.34, 231.65) = 35.54$, $p < .001$, $\eta^2 = .26$. In general, *RMSE* increased as person

521  location moved towards either tail of the distribution, away from the average IQ of 100. This

522  effect, also seen in the other simulated populations, is visualized as a parabolic shape in

523  Figure 4. By contrast, in the analysis with *MSD*, the main effect of person location is

524  visualized as a sinusoidal pattern (see Figure 5 and discussion section below). This effect of

525  person location on norming bias is a previously reported feature of continuous norming

526  procedures (cf. A. Lenhard et al., 2019). As such, this effect is not directly relevant to the

527  question of whether weighting, per se, reduces norm bias due to non-representative sampling.

528  What is important to note (and is readily seen in Figures 4 and 5) is that WCN and SCN

529  perform equally well, in terms of error measures, when processing normative samples drawn

530  from a demographically representative, reference population. This makes intuitive sense,

531  because with representative samples, there are no cell-wise departures from expected

532  demographic proportions, to which weights could be applied to correct for bias in the

533  norming process.

534  *Population 2: Mild under-representation of high education*

535       In samples drawn from Population 2, we found a main effect of norming method,

536  *RMSE*: $F(1, 99) = 74.94$, $p < .001$, $\eta^2 = .43$, *MSD*: $F(1, 99) = 1924.64$, $p < .001$, $\eta^2 = .95$.

537  WCN was superior to SCN in reducing norm bias resulting from non-representative samples.

538  With *RMSE*, we also observed an interaction between person location and norming method,

539  $F(2.72, 268.76) = 41.72$, $p < .001$, $\eta^2 = .30$. As shown in Figure 4, WCN reduced the error

540  measure to a greater degree in the upper range of person location than in the lower range. In

541  Population 2, individuals of higher education (and consequently, higher cognitive ability) are

542  under-represented. Thus, the interaction shows that WCN is correcting for norm bias in the

543  region of person location that is under-represented in the normative samples.

544          In the comparison of WCN in Population 2 to the benchmark of SCN in Population 1,

545  there was no main effect of population on *RMSE*. That is, even under the conditions of non-

546  representativeness in Population 2, WCN did not differ from the benchmark on the error

547  measure. This suggests that weighting successfully compensated for any norm bias due to

548  demographic non-representativeness in Population 2, when that bias was measured by *RMSE*.

549  The results differed for *MSD,* where we observed a main effect of population, $F(1, 198) =$

550  15,37, $p < .001$, $\eta^2 = .07$, and an interaction between population and person location, $F(2.31,$

551  $456.93) = 2.91$, $p = .048$, $\eta^2 = .01$. These findings indicated that, with respect to *MSD*, WCN

552  did not fully correct norm bias in samples from Population 2. In addition, this difference in

553  MSD between WCN and the benchmark condition[9] exceeded the threshold of practical

554  relevance in the lower region of person location, with absolute differences of 0.68 IQ points

555  at IQ 70 and 0.99 IQ points at IQ 62.5.

556  ***Population 3: Moderate under-representation of high education***

557          For both error measures, the ANOVAs with normative samples drawn from

558  Population 3 yielded significant main effects of norming method, *RMSE*: $F(1, 99) = 155.76$,

559  $p < .001$, $\eta^2 = .61$, *MSD*: $F(1, 99) = 2707.76$, $p < .001$, $\eta^2 = .97$, and significant interactions

560  between norming method and person location, *RMSE*: $F(2.79, 276.13) = 71.89$, $p < .001$, $\eta^2 =$

561  .42, *MSD*: $F(2.63, 260.05) = 6.58$, $p = .001$, $\eta^2 = .06$. These analyses produced larger effect

562  sizes than those for Population 2, mirroring the difference in representativeness between the

563  two populations. This suggests that WCN exerts a larger corrective effect on norm-score bias

564  with normative samples that display greater deviations from demographic representativeness.

---

[9] In the remainder of this section, the absolute difference in error measure between
norming methods will be labeled |ΔRMSE| or |ΔMSD|, as appropriate.

565       In comparing WCN in Population 3 to SCN in Population 1, we observed significant

566    main effects of population, $RMSE$: $F(1, 198) = 9.52$, $p = .002$, $\eta^2 = .05$, $MSD$: $F(1, 198) =$

567    12.52, $p = .001$, $\eta^2 = .06$, and significant interactions between population and person location,

568    $RMSE$: $F(2.81, 557.01) = 3.06$, $p = .031$, $\eta^2 = .02$, $MSD$: $F(2.26, 447.76) = 3.28$, $p = .033$, $\eta^2$

569    $= .02$. In the normative samples drawn from Population 3, WCN yielded greater norming

570    error than the benchmark within the low and high regions of person location. However, a

571    practically significant difference in IQ scores occurred at only one location for $RMSE$ (IQ

572    137.5, $|\Delta RMSE| = 0.57$ IQ points), and one location for $MSD$ (IQ 62.5, $|\Delta MSD| = 0.68$ IQ

573    points).

574    *Population 4: Under-representation of both low and high education*

575       As with Populations 2 and 3, the analyses of normative samples drawn from

576    Population 4 produced significant main effects of norming method, $RMSE$: $F(1, 99) = 39.02$,

577    $p < .001$, $\eta^2 = .28$, $MSD$: $F(1, 99) = 164.63$, $p < .001$, $\eta^2 = .62$, and significant interactions

578    between norming method and person location, $RMSE$: $F(3.02, 299.01) = 42.43$, $p < .001$, $\eta^2 =$

579    .30, $MSD$: $F(2.53, 250.47) = 42.44$, $p = .001$, $\eta^2 = .30$. However, with Population 4, where

580    both tails of the education distribution were under-represented, WCN did not provide greater

581    reduction of norm-score bias than SCN. The interactions revealed that in the low region of

582    person location $RMSE$ was greater for WCN than SCN (IQ 77.5 $|\Delta RMSE| = 0.61$ IQ points,

583    IQ 70.0 $|\Delta RMSE| = 1.37$ IQ points, IQ 62.5 ($|\Delta RMSE| = 1.43$ IQ points). For $MSD,$ the

584    interaction between norming method and person location was more complex, with SCN

585    providing greater reduction of norm bias than WCN in the low region of person location (IQ

586    70.0: $|\Delta MSD| = 1.25$ IQ points; IQ 62.5: $|\Delta MSD| = 2.00$ IQ points), and WCN outperforming

587    SCN in the high region of person location (IQ 115.0: $|\Delta MSD| = 0.60$ IQ points; IQ 122.5:

588    $|\Delta MSD| = 0.90$ IQ points; IQ 130.0: $|\Delta MSD| = 0.87$ IQ points; IQ 137.5: $|\Delta MSD| = 0.88$ IQ

589    points).

590        In comparing WCN in Population 4 to SCN in Population 1, we found significant

591 main effects of population, *RMSE*: $F(1, 198) = 38.46$, $p < .001$, $\eta^2 = .16$, *MSD*: $F(1, 198) =$

592 $31.10$, $p < .001$, $\eta^2 = .14$, and significant interactions between population and person location,

593 *RMSE*: $F(2.83, 560.75) = 22.19$, $p < .001$, $\eta^2 = .10$, *MSD*: $F(2.55, 505.02) = 19.44$, $p < .001$,

594 $\eta^2 = .09$. For *MSD*, we examined the simple effects underlying the interaction and found that

595 WCN reduced norm bias more than the benchmark only at the highest person locations (IQ

596 >122.5). At IQ 122.5, the difference was 0.58 IQ points.

597 *Populations 5 (Biased joint distributions) and 6 (Clustered distributions)*

598        For populations 5 and 6, the ANOVAs revealed no main effects of norming method.

599 With respect to *RMSE*, the differences between WCN and SCN did not exceed 0.5 IQ points

600 at any point in the range of person location. When we compared WCN in Population 5 to the

601 benchmark (SCN in Population 1), the ANOVA for *MSD* returned a significant main effect of

602 population, $F(1, 198) = 5.02$, $p = .026$, $\eta^2 = .03$, and a significant interaction between

603 population and person location, $F(2.29, 453.98) = 6.25$, $p = .001$, $\eta^2 = .03$. At some person

604 locations, the absolute difference between WCN and the benchmark exceeded 0.5 IQ points.

605 However, the valence of these differences varied over the range of person location. At IQ

606 122.5 and IQ 130.0, MSD was closer to zero for WCN than for the benchmark, but at IQ

607 137.5 this pattern was reversed ($\Delta MSD| = 0.73$). The results suggest that weighting offers no

608 clearcut advantage in reducing the error associated with norming, when normative samples

609 are drawn from a population with biased joint distributions of the demographic variables.

610                                   **Discussion**

611 **Summary of results**

612        The present study examined whether compensatory weighting at the raw score level,

613 when combined with continuous norming procedures, would reduce bias in norm scores

614 derived from demographically non-representative norm samples. To pursue this aim, we

615     simulated six populations in which the distributions of demographic variables departed to

616     various degrees from expected proportions. We modeled a latent cognitive ability, which we

617     used as the input for a one-parameter logistic IRT model to create raw test scores. We drew

618     normative samples from the six populations, and generated IQ-type norm scores by applying

619     weighted continuous norming (WCN) and semi-parametric continuous norming without

620     weighting (SCN). We used mean square error (*RMSE*) and mean signed difference (*MSD*) as

621     measures of norm-score bias.

622         Our first hypothesis proposed that when processing non-representative normative

623     samples, WCN would produce less-biased norm scores than SCN. The predicted advantage of

624     WCN was most apparent in samples drawn from populations 2 and 3, in which individuals

625     with high levels of education were under-represented. In samples drawn from populations 4

626     (Under-representation of both low and high education) and 6 (Clustered sampling), WCN

627     showed no benefit over SCN, but neither did it degrade the quality of norm scores, relative to

628     continuous norming without compensatory weighting. In population 5 (Biased joint

629     probabilities), we found that WCN led to a small increase in norm score error, but only at

630     certain points in the range of cognitive ability.

631         In normative samples drawn from Population 1, which served as the standard of

632     representativeness for the demographic variables, WCN demonstrated no advantage over

633     SCN. This result is not surprising: WCN creates weights to compensate for departures from

634     representativeness.  Because Population 1 was the benchmark, in terms of demographic

635     composition, random samples drawn from it were expected to be demographically

636     representative.

637         Population 2 introduced deviations from the benchmark distribution of education, the

638     demographic variable with the strongest effect on cognitive ability. Specifically, level 1 (high

639     education/high ability) was mildly under-represented, and level 3 (low education/low ability)

640  was proportionately over-represented. In samples drawn from Population 2, WCN yielded

641  greater reduction in norm-score bias than SCN, for both error measures, across the entire

642  range of cognitive ability.

643      Population 3 presented a pattern of non-representativeness on education that was

644  similar to that in Population 2, but greater in magnitude. The comparison of normative

645  samples drawn from Populations 2 and 3 was relevant to testing our second hypothesis,

646  which specified that as the non-representativeness of the normative sample increased, norm-

647  score bias would increase for both methods, but that the increase in bias would be smaller for

648  WCN than for SCN. Our findings provided support for this hypothesis: with samples drawn

649  from Population 3, the magnitude of norming error for WCN was larger than it was in the

650  Population 2 analyses, although WCN retained its superiority to SCN in terms of reducing

651  norm score bias. With population 3, the increase in norming error associated with WCN

652  depended on person location – it occurred at either extreme of the range of the cognitive

653  ability variable, but not in the middle region. In no instances, however, did these increases in

654  the error measures exceed 1 IQ point.

655      Population 4 embodied a further scenario of demographic non-representativeness, in

656  which both tails of the education distribution were under-represented, and the central region

657  of the distribution was proportionately over-represented. In terms of the average degree of

658  misrepresentation across the three levels of education, Population 4 did not differ from

659  Population 3. Where the effect of the demographic manipulation differs is on the raw score

660  distributions. In Population 4, the manipulation attenuates the variance of the raw score

661  distributions, because under-sampling both tails of the education distribution results in an

662  under-sampling of the very high and low raw scores that reside in those regions. In addition,

663  whereas in Population 3 the pattern of misrepresentation affects the mean of the raw score

664  distribution, in Population 4 the mean is not affected, because there is equal under-

665  representation of both tails of the raw score distribution.

666      In normative samples drawn from Population 4, we observed that in certain regions of

667  person location, WCN was less effective than SCN in reducing norm-score bias, which is

668  consistent with our third hypothesis. Specifically, we found that the disparity between WCN

669  and SCN increased at both tails of the cognitive ability distribution, with WCN showing the

670  greatest magnitude of norming error in the lowest region of person location.

671      To put this finding into context, consider how the raw score distribution of a

672  demographic subgroup is affected differentially by adding additional individuals, as opposed

673  to weighting the existing raw scores without increasing sample size. Adding more individuals

674  increases the variance of the raw score distribution, whereas weighting existing raw scores

675  does not affect the variance. In Population 4, furthermore, the variance of the low and high

676  ability groups was reduced by the pattern of under-representation, which results in fewer

677  individuals in each of these groups. Therefore, weighting the raw scores of the under-

678  represented groups increases the influence of any sampling error that exists in the raw score

679  distributions. This phenomenon may explain our finding that WCN resulted in greater

680  norming error in the under-represented, low region of person location. By contrast, WCN did

681  not yield increased norming error in the central region of person location, where there are

682  more observations present and a consequent reduction in sampling error. Our findings

683  suggest that researchers should employ WCN with caution when processing normative

684  samples where the non-representative subgroups are also those containing few individuals.

685      In Population 5, the joint distributions of the demographic variables (resulting from a

686  complete cross-classification of the three variables) were manipulated in a pattern of

687  alternating over- and under-representation. This was accomplished so that the marginal

688  distributions of the variables closely approximated those of the reference population. Thus,

36

689     Population 5 simulates a sampling scenario wherein demographic misrepresentation occurs at

690     a level that is "beyond the reach" of cNORM's raking and weighting methods, which operate

691     only on marginal distributions.

692         Under these conditions, WCN did not provide any improvement in the reduction of

693     norm-score bias over SCN. However, our manipulation of the joint probabilities, as it turned

694     out, did not strongly affect the means and variances of the raw score distributions. Thus, this

695     analysis leaves unanswered the question of how WCN might perform when misrepresentation

696     at the level of joint probabilities does bias the parameters of the raw score distributions.

697         It is important to keep in mind that in our simulation study, the three demographic

698     variables were modeled so that education had the strongest relationship with cognitive ability,

699     and thus had more impact on norm score accuracy than ethnicity or region. Thus, our findings

700     with Population 5 do not reflect the range of possible relationships between demographic

701     factors and cognitive ability (e.g., other variables that are highly correlated with ability, or

702     that interact with each other). In these alternate scenarios, it is unknown how

703     misrepresentation in the joint distributions might affect raw score means and variances. Later

704     in this section, we provide guidance on how to address these scenarios in practice.

705         In Population 6 (clustered distributions), the distributions of the demographic

706     variables were manipulated *within* each of the six age cohorts. This manipulation is best

707     understood in comparison to Population 1, in which the marginal and joint probabilities of the

708     entire population are replicated within each age cohort. In Population 6, by contrast, two-

709     thirds of the joint distribution cells contained no data, meaning that the overall demographic

710     distributions were *not* replicated *within* the age cohorts. However, the pattern of data deletion

711     was such that the marginal and joint probabilities of Population 6, averaged over the entirety

712     of the population (across all age cohorts), matched those of Population 1.

**Commented [DH22]:** I found the original line of argument about the Population 5 findings to be convoluted and difficult to follow. I tried to simplify the argument here, but please check against the original to make sure that the new narrative accurately reflects the points you were trying to make, and that it omits no important details.

713    In normative samples drawn from Population 6, the age-specific patterns of

714    demographic non-representativeness affected the parameters of the raw score distributions

715    within each age cohort. Raking per se cannot compensate for discrepancies of this nature,

716    because raking operates on marginal probabilities of the entire normative sample, not those

717    within each age cohort. It is therefore counter-intuitive to find, as we did, that neither WCN

718    nor SCN yielded increases in norm-score bias, when compared to the benchmark condition.

719    We attribute this finding to the influence of the semi-parametric continuous norming method

720    that underlies both WCN and SCN. As noted previously, this method models a developing

721    cognitive ability as a monotonic function of age and person location. It thereby uses the stable

722    variance of the entire normative sample to smooth the parameters of the raw score

723    distributions across age cohorts, even when those age-specific distributions are affected by

724    varying levels of demographic non-representativeness.

**Commented [DH23]:** Please compare this paragraph to the original manuscript. The original used jargon such as "wavy iso-percentiles" and "stiffness of the method" which was unfamiliar to me. I tried to rewrite it in a simpler fashion that relied on basic concepts of continuous norming that were introduced earlier in the manuscript. Please confirm that I'm capturing your intended meaning here.

725    **Implications for the use of WCN in test norming**

726    Our study showed that WCN reduces norm-score bias under certain patterns of non-

727    representativeness of a demographic variable, where that variable is strongly correlated with

728    the test score being normed. The pattern of results across the six simulated populations,

729    however, suggested that even when a demographic variable has a strong effect on raw scores,

730    it produces relatively small distortions in resulting norm scores. Even under conditions

731    representing large departures from demographic representativeness, the differences in *RMSE*

732    between norm scores derived using WCN and those from representative samples did not

733    exceed 2 IQ points.

734    With norming methods that generate norms independently for each age group, we

735    would expect departures from demographic representativeness to cause greater levels of

736    norm-score bias. These conventional norming methods lack the previously noted advantage

737    of continuous norming, which can smooth out local effects of non-representativeness.

738      Consistent with this view, we have demonstrated previously that with conventional norming

739      per age group, *RMSE* is about twice as high, on average, as with semi-parametric continuous

740      norming, even with representative random samples (W. Lenhard & Lenhard, 2021). The

741      selection of an appropriate norming method is therefore a critical prerequisite for accurate

742      test norms, regardless of whether this procedure is used with or without weighting. By

743      contrast, the size of the normative sample is less critical, if continuous norming is used. For

744      example, we found that increasing sample size from 100 to 250 per age group did not yield

745      significant reduction in *RMSE*, when continuous norming methods were used (A. Lenhard et

746      al., 2019).

747          Clearly, the best practice is to prevent problems associated with non-

748      representativeness in the first place, by collecting an adequately sized, demographically

749      representative sample for norming. Post-hoc weighting procedures are no substitute for a

750      well-planned data collection effort that draws randomly from the general population. Care

751      must also be taken to avoid over-sampling from clinical settings, as this will bias the sample

752      towards individuals of lower ability.

753          The current study demonstrates the utility of weighting procedures in reducing norm-

754      score error under conditions of mild-to-moderate non-representativeness of a demographic

755      variable. Nevertheless, we also found that the ability of WCN to reduce norm-score bias was

756      degraded, when we reduced the marginal probability of the high level of education to 20%

757      from the reference value of 40% (that is, when the size of that subgroup was half that needed

758      for a representative sample). Our work further shows that the effectiveness of weighting

759      depends on the location of under-represented demographic groups on the spectrum of person

760      ability. With a typical cognitive ability that is normally distributed in the general population,

761      random sampling will yield relatively small subgroups at either tail of the ability distribution.

762      If these extreme subgroups are under-sampled to begin with, any sampling error embodied in

763    the raw score distributions will only be multiplied by the application of compensatory

764    weights. This can lead to increased norm-score bias, as illustrated in our results. The remedy,

765    of course, is to ensure that these low- and high-ability groups are represented in adequate

766    numbers.

767          As described previously, the raking procedures used in this study operate only on the

768    marginal distributions of the demographic variables. Census information on the joint

769    distributions of the demographic variables (e.g., the expected probability for the joint

770    category of low education/non-white ethnicity) is not always available. However, when that

771    information is available, it can be incorporated in the raking procedures through a recoding

772    process. For example, the crossing of two demographic variables, each of which has three

773    categories, results in nine cross-classification cells. These classifications can be recoded into

774    nine levels of single dummy variable. The expected joint probabilities of the cross-classified

775    cells thereby become the expected marginal probabilities of the dummy variable. The risk in

776    this approach comes from increasing the number of categories, which also increases the

777    likelihood that one or more category would have a very low expected probability. Under

778    these circumstances, of course, even adequately sampled categories may hold only a few

779    individuals, thus increasing the influence of sampling error when weights are applied.

780          To counter this tendency, we often recommend reducing the number of demographic

781    categories by combining groups that are not expected to differ significantly in mean location

782    on the ability variable. This practice can be applied to either marginal categories or joint

783    cross-classifications, when the latter are subject to the recoding procedure described in the

784    previous paragraph.

785    **Limitations of the study**

786          This study evaluated only one method of post-stratification: raking with marginal

787    probabilities as the input. We did not examine fully cross-classified post-stratification (i.e., a

788　method that takes joint probabilities into account). Instead, we analyzed norm samples drawn

789　from Population 5 (biased joint distributions), to determine the performance of raking under

790　conditions where the marginal probabilities are representative, but the joint probabilities are

791　not. In Population 5, we did not find that WCN, which includes raking, yielded increased

792　norm-score bias compared to the benchmark condition. This may have been due to the

793　magnitude of non-representativeness in the cross-classification cells. The demographic

794　deficiencies in these cells may not have been great enough to expose the inability of raking to

795　compensate for such deficiencies.

796　　　　In our study, we simulated three demographic variables (education, ethnicity, region),

797　with varying levels of correlation with the latent cognitive ability (strong, moderate, weak,

798　respectively). We did not model any interactions among these three variables. Demographic

799　variables that interact in their effects on cognitive ability might yield larger disturbances in

800　the raw score distributions of the cross-classification cells. Under these conditions, as we

801　have demonstrated, weighting carries the risk of increasing norm-score bias. However, the

802　main effect of education on test scores in our study probably represents the upper limit of

803　analogous effects that could occur in real-world normative samples. With demographic

804　variables that have smaller effect sizes, of course, we can expect resulting norm-score biases

805　to also diminish in magnitude.

806　　　　A second limitation was that our study modeled only one latent psychological

807　variable: a cognitive ability that increases monotonically with increasing age. Other variables

808　measured by psychometric tests (e.g., the "big five" personality traits, Donnellan & Lucas,

809　2008) may not manifest the same dependency on age, and they may be affected by

810　demographic variables with different characteristics than the ones simulated in our study.

811　When norming tests of personality traits, therefore, it may be appropriate to apply a

812　weighting method that is not combined with continuous norming procedures.

813    A third caution relates to the mathematical underpinnings of the cNORM norming

814    process. cNORM uses a semi-parametic continuous norming method that requires the

815    expansion of a Taylor polynomial (for more details, see [INSERT CITATION]). The

816    modeling process calls for specification of a parameter ($k$), that sets an upper bound on the

817    exponents of person location and age. In the current study, we used a default value of $k = 4$

818    for both location and age. It is possible that more precise models of the latent cognitive

819    ability could have been obtained with different values of $k$. Simulation studies published

820    elsewhere ([INSERT CITATION]) have compared norm-score bias across a range of values

821    of $k$. These findings suggest that $k = 5$ for location and $k = 3$ for age provide an optimal

822    balance between norm score accuracy and processing load. As a result, we have selected

823    these values as the defaults for the current version of cNORM.

824    Finally, our study examined weighting only as applied to the semi-parametric

825    continuous norming method implemented in the cNORM package. We did not combine

826    weighting with other continuous norming approaches, such as parametric continuous norming

827    (e.g., Stasinopoulos et al., 2018).  Earlier in this paper, we pointed out that the semi-

828    parametric continuous method, because it does not rely on splines to model age-related

829    changes in ability, may be better suited for certain conditions of non-representativeness in

830    normative samples (e.g., the clustered distributions modeled in Population 6). Moreover, we

831    have demonstrated elsewhere (see A. Lenhard, 2019) that the cNORM approach yields less

832    norm-score bias than parametric continuous norming with skewed raw score distributions,

833    and with sample sizes of 150 or less per age group. Yet, the efficiency of post-stratification

834    techniques combined with parametric continuous norming remains to be investigated.

835    **Concluding Remarks and Outlook**

836    The application of weighting techniques to the norming of psychometric tests is a

837    relatively new area of study. Unsurprisingly, therefore, several additional research questions

838     emerged from the current simulation protocol. For example, we implemented raking weights

839     twice within cNORM: Once during ranking of raw scores, and then again during regression

840     modeling. But we did not evaluate the relative value, in terms of reducing norm-score error,

841     of the second step. It is therefore possible that applying weights to the regression analysis was

842     of little benefit, or that it may have even increased norm score bias. The latter might occur

843     because, as noted previously, weighting can multiply the effects of sampling error in under-

844     represented groups.

845         This last caution serves as a final reminder that weighting techniques are no substitute

846     for the painstaking process of assembling a demographically representative normative

847     sample. Our study has shown, however, that if such samples still exhibit reasonably small

848     departures from representativeness, the weighting methods implemented in cNORM offer a

849     useful way of mitigating any resulting norm-score bias.

853 **References**

854 Cole T. (1988). Fitting smoothed centile curves to Reference Data. *Journal of the Royal*

855 *Statistical Society Series A (Statistics in Society), 151*(3), 385.

856 Cole, T. J., & Green P. J. (1992). Smoothing reference centile curves: The lms method and

857 penalized likelihood. *Statistics in Medicine, 11*(10), 1305–19.

858 Donnellan, M. B., & Lucas, R. E. (2008). Age differences in the Big Five across the life span:

859 evidence from two national samples. *Psychology and aging, 23*(3), 558–566.

860 https://doi.org/10.1037/a0012897

861 Ireland, C. T., & Kullback, S. (1968). Contingency tables with given marginals. *Biometrika,*

862 *55*(1), 179–188. https://doi.org/10.1093/biomet/55.1.179

863 Kalton, G., & Flores-Cervantes, I. (2003). Weighting Methods. *Journal of Official Statistics,*

864 *19*(2), 81-97.

865 Kaufman, A. S., & Kaufman, N. L. (2004). *Kaufman Assessment Battery for Children Second*

866 *Edition.* San Antonio: Pearson Clinical Assessment.

867 Kennedy, L., & Gelman, A. (2021). Know your population and know your model: Using

868 model-based regression and poststratification to generalize findings beyond the

869 observed sample. *Psychological Methods, 26*(5), 547–558.

870 https://doi.org/10.1037/met0000362

871 Lenhard, A. & Lenhard, W. & Gary, S. (2018). cNORM. (Continuous Norming). *The*

872 *Comprehensive R Network.* https://cran.r-

873 project.org/web/packages/cNORM/index.html. doi:10.1177/1073191116656437

874 Lenhard, A., Lenhard, W. & Gary, S. (2019). Continuous norming of psychometric tests: A

875 simulation study of parametric and semi-parametric approaches. *PLOS ONE, 14*(9),

876 e0222279. https://doi.org/10.1371/journal.pone.0222279

877    Lenhard, A., Lenhard, W., Suggate, S. & Segerer, R. (2018). A continuous solution to the

878        norming problem. *Assessment, 25*(1), 112-125. doi.org/10.1177/1073191116656437

879    Lenhard, A., Lenhard, W., Segerer, R. & Suggate, S. (2015). *Peabody Picture Vocabulary*

880        *Test - Revision 4 (PPVT-4)*, German adaptation. Pearson Assessment.

881    Lenhard, W. & Lenhard, A. (2021). Improvement of Norm Score Quality via Regression-

882        Based Continuous Norming. *Educational and psychological measurement, 81*(2),

883        229-261. https://doi.org/10.1177/0013164420928457

884    Lenhard, W., Lenhard, A. & Schneider, W. (2017). *ELFE II - Ein Leseverständnistest für*

885        *Erst- bis Siebtklässler.* nHogrefe.

886    Little, R. J. (1993). Post-stratification: A modeler's perspective. *Journal of the American*

887        *Statistical Association, 88*(423), 1001-1012. https://doi.org/10.2307/2290792

888    Lumley, T. (2011). *Complex Surveys – A Guide to Analyses Using R*. Wiley.

889    Lumley, T. (2017). leaps: Regression Subset Selection. *The Comprehensive R Network.*

890        https://cran.r-project.org/web/packages/leaps/index.html

891    Park, D. K., Gelman, A., & Bafumi, J. (2004). Bayesian multilevel estimation with post-

892        stratification: State-level estimates from national polls. *Political Analysis, 12*(4), 375-

893        385.

894    Price-Mohr, R., Price, C. (2017). Gender Differences in Early Reading Strategies: A

895        Comparison of Synthetic Phonics Only with a Mixed Approach to Teaching Reading

896        to 4–5 Year-Old Children. *Early Childhood Education Journal, 45*, 613–620.

897        https://doi.org/10.1007/s10643-016-0813-y

898    Sang-Wook Y., Heechoul O., Soon-Ae S., & Jee-Jeon Y. (2015). Sex-age-specific

899        association of body mass index with all-cause mortality among 12.8 million Korean

900        adults: a prospective cohort study, *International Journal of Epidemiology, 44*(5),

901        1696–1705. https://doi.org/10.1093/ije/dyv138

902  Stasinopoulos M. D., Rigby, R. A., Voudouris, V., Akantziliotou, C., Enea, M. & Kiose, D.

903      (2018). gamlss: Generalised Additive Models for Location Scale and Shape. *The*

904      *Comprehensive R Network.* https://cran.r-project.org/web/packages/gamlss/index.html

905  Wechsler, D. (1939). *The measurement of adult intelligence.* Williams & Wilkins Co.

906      https://doi.org/10.1037/10020-000

907  Wechsler, D. (2008). *WAIS-IV Technical and interpretive manual*. Pearson.

908  Wechsler, D. (2014). *WISC-V Technical and interpretive manual.* Pearson.

909  Zachary, R. A., & Gorsuch, R. L. (1985). Continuous norming: Implications for the WAIS-R.

910      *Journal of Clinical Psychology, 41*(1), 86–94.

**Figure Captions**

911

912      *Figure 1.* Modeled cognitive ability in reference population.

913      *Figure 2.* RMSE across simulated populations, with (WCN) or without (SCN)

914 weighting. The grey rectangles represent 95% confidence intervals.

915      *Figure 3.* MSD across simulated populations with (WCN) or without (SCN)

916 weighting. The grey rectangles represent 95% confidence intervals.

917      *Figure 4.* RMSE across simulated populations, with (WCN) or without (SCN)

918 weighting, as a function of person location. The dotted grey line represents SCN with norm

919 samples drawn from Population 1 (benchmark).

920      *Figure 5.* MSD across simulated populations, with (WCN) or without (SCN)

921 weighting, as a function of person location. The dotted grey line represents SCN with norm

922 samples drawn from Population 1 (benchmark).