

Using Heteroskedastic Ordered Probit Models to Recover Moments of Continuous Test Score Distributions from Coarsened Data

AUTHORS

Sean F. Reardon

Stanford University

Benjamin R. Shear

Stanford University

**Katherine E.
Castellano**

Educational Testing Service

Andrew D. Ho

Harvard Graduate School of
Education

ABSTRACT

Test score distributions of schools or demographic groups are often summarized by frequencies of students scoring in a small number of ordered proficiency categories. We show that heteroskedastic ordered probit (HETOP) models can be used to estimate means and standard deviations of multiple groups' test score distributions from such data. Because the scale of HETOP estimates is indeterminate up to a linear transformation, we develop formulas for converting the HETOP parameter estimates and their standard errors to a scale in which the population distribution of scores is standardized. We demonstrate and evaluate this novel application of the HETOP model with a simulation study and using real test score data from two sources. We find that the HETOP model produces unbiased estimates of group means and standard deviations, except when group sample sizes are small. In such cases, we demonstrate that a "partially heteroskedastic" ordered probit (PHOP) model can produce estimates with a smaller root mean squared error than the fully heteroskedastic model.

VERSION

June 2016

Suggested citation: Reardon, S.F., Shear, B.R., Castellano, K.E., & Ho, A.D. (2016). Using Heteroskedastic Ordered Probit Models to Recover Moments of Continuous Test Score Distributions from Coarsened Data (CEPA Working Paper No.16-02). Retrieved from Stanford Center for Education Policy Analysis: <http://cepa.stanford.edu/wp16-02>

**Using Heteroskedastic Ordered Probit Models
to Recover Moments of Continuous Test Score Distributions from Coarsened Data**

Sean F. Reardon (Stanford University)

Benjamin R. Shear (Stanford University)

Katherine E. Castellano (Educational Testing Service)

Andrew D. Ho (Harvard Graduate School of Education)

DRAFT: June 2016

Forthcoming, *Journal of Educational and Behavioral Statistics*

Direct correspondence to Sean F. Reardon (sean.reardon@stanford.edu). This research was supported by grants from the Institute of Education Sciences (#R305D110018 and #R305B090016) and a grant from the Spencer Foundation (#201500058). The paper benefited from collaboration with Erin Fahle and thoughtful comments from J.R. Lockwood. We thank Demetra Kalogrides for excellent research assistance, and Richard Williams for his responsiveness and assistance regarding our inquiries about his `og1m` Stata package. Some of the data used in this paper were provided by the National Center for Education Statistics (NCES). The opinions expressed here are ours and do not represent views of NCES, the Institute of Education Sciences, the Spencer Foundation, or the U.S. Department of Education.

Using Heteroskedastic Ordered Probit Models to Recover Moments of Continuous Test Score Distributions from Coarsened Data

Abstract

Test score distributions of schools or demographic groups are often summarized by frequencies of students scoring in a small number of ordered proficiency categories. We show that heteroskedastic ordered probit (HETOP) models can be used to estimate means and standard deviations of multiple groups' test score distributions from such data. Because the scale of HETOP estimates is indeterminate up to a linear transformation, we develop formulas for converting the HETOP parameter estimates and their standard errors to a scale in which the population distribution of scores is standardized. We demonstrate and evaluate this novel application of the HETOP model with a simulation study and using real test score data from two sources. We find that the HETOP model produces unbiased estimates of group means and standard deviations, except when group sample sizes are small. In such cases, we demonstrate that a "partially heteroskedastic" ordered probit (PHOP) model can produce estimates with a smaller root mean squared error than the fully heteroskedastic model.

Keywords: Heteroskedastic ordered probit model; test score distributions; coarsened data.

Using Heteroskedastic Ordered Probit Models to Recover Moments of Continuous Test Score Distributions from Coarsened Data

The widespread availability of aggregate student achievement data provides a potentially valuable resource for researchers and policymakers alike. Often, however, these data are only publicly available in “coarsened” form in which students are classified into one or more ordered “proficiency” categories (e.g., “Basic,” “Proficient,” “Advanced”). Although proficiency category data are clearly useful when proficiency status itself is of substantive interest, coarsened data pose challenges for the analyst when moments of the underlying test score distributions are of interest. Proficiency rates convey information about a single point in a cumulative test score distribution. This not only limits the information available to the analyst about the underlying test score distribution, but it also complicates inferences about relative changes in achievement levels in different population subgroups, a point illustrated by Ho (2008) and Holland (2002).

For example, suppose one wants to compare the average test scores among multiple schools, but one knows only the proportion of students scoring in each of several ordered proficiency categories. If the underlying test score distributions have unequal variances among schools, then rankings of schools on the basis of the percentages scoring at or above a given proficiency category will depend on which threshold is chosen. Moreover, rankings of schools based on percentages above some threshold will not, in general, match rankings based on mean scores. The same problem holds if one wishes to compare average test scores among multiple student subgroups (such as racial/ethnic groups) or to compare average test scores in a given school over time. In each case, judgements about the relative magnitude of between-group differences and even the ordering of groups’ average performance will be dependent on what proficiency category threshold is used. These and other limitations posed by the coarsening of standardized test scores have been described extensively (Ho, 2008; Ho & Reardon, 2012; Holland, 2002;

Jacob, Goddard, & Kim, 2013; Jennings, 2011).

With access to only coarsened test score data, therefore, comparisons of average performance among groups of students may be ambiguous. Unfortunately, most publicly available data on student performance on state standardized tests consists of coarsened test scores. Most states, for example, do not report school- or district-level test score means (and very few report standard deviations). The *EDFacts* Assessment Database (e.g., U.S. Department of Education, 2015), for example, provides test score data for every public school in the United States, but does not include means and standard deviations. Rather, it contains the counts of students (by school, grade, subject, and student subgroup) scoring in each of two to five state-defined performance levels, as required under the *Elementary and Secondary Education Act* (ESEA). While these data are a valuable resource for educators, policymakers, and researchers, their utility is severely hampered by the absence of test score means and standard deviations.

In this paper, we describe an approach that allows the analyst to recover more complete information about continuous test score distributions when only coarsened test score data are available. To achieve this, we propose a novel application of the heteroskedastic ordered probit (HETOP) model (e.g., Alvarez & Brehm, 1995; Greene & Hensher, 2010; Keele & Park, 2006; Williams, 2009). As we describe, the HETOP model can be used to recover means and standard deviations of continuous test score distributions of multiple groups from coarsened data. These groups may be schools, districts, or demographic subgroups. Estimates of these group means and standard deviations can be used to estimate intraclass correlations (ICCs), between-group achievement gaps, and other theoretically interesting or policy-relevant statistics, just as if each group's mean and standard deviation were provided directly.

The methods we describe generalize prior work quantifying achievement gaps in an ordinal or nonparametric framework, both with continuous (Ho, 2009) and coarsened (Ho & Reardon, 2012;

Reardon & Ho, 2015) test scores. Although we describe the use of such models to recover moments of test score distributions from aggregate proficiency data, the methods are applicable to other educational testing contexts when only coarsened scores are reported, such as Advanced Placement (AP) or English language proficiency exams. Aggregate data on coarse ordered scales can also arise in college rankings, health research and practice scales, and income reporting. In all of these cases, our methods enable the estimation of group means and standard deviations from ordered data.

The paper is organized into four main sections. In Section 1, we describe the statistical and conceptual framework for our application of the HETOP model. In Section 2, we use Monte Carlo simulations to evaluate recovery of the parameters of interest across a range of scenarios that might be encountered in practice. In Section 3, we use two real test score data sets, one from the National Assessment of Educational Progress (NAEP) and one from a State testing program, to evaluate the extent to which the key assumption of the HETOP model holds for real data. For these case studies, both student-level scale scores and coarsened proficiency counts are available, allowing us to evaluate the agreement between HETOP model estimates of means and standard deviations and estimates of the same parameters based on uncoarsened scale score data. Section 4 summarizes and discusses the results and offers recommendations for applying the methodology in practice.

1 Background and Statistical Framework

1.1 Canonical Application: Data, Assumptions, and Estimands

In our context of interest—the reporting of large-scale educational test scores in proficiency categories—the data consist of frequencies of students scoring within each of K ordered categories (often called performance levels) for G groups. Groups might be defined by racial/ethnic categories, schools, school districts, or other relevant categories. Such data can be summarized in a $G \times K$ matrix. Each cell in the matrix indicates the observed frequency of students from group $g = \{1, \dots, G\}$ scoring at

performance level $k = \{1, \dots, K\}$ of a test. The performance levels describe ordered degrees of student proficiency in a content area. In standard current state testing practice, a panel of content experts selects one to four cutscores that divide the score scale into performance levels through a standard setting procedure (e.g., Cizek, 2012).

Let y denote a continuous random variable (scaled test scores, in our example), with μ_g and σ_g denoting the mean and standard deviation, respectively, of y in group g . Although we make no specific distributional assumptions about the shape of the distributions of y in each group, we do make the assumption that the distributions are “respectively normal” (Ho, 2009; Ho & Reardon, 2012). This means we assume there exists a continuous monotonic increasing function f defined for all values of y , such that $y^* = f(y)$ has a normal distribution within each group g :

$$y^* | (G = g) \sim N(\mu_g^*, \sigma_g^{*2}) \quad (1)$$

This does not require that the test score y be normally distributed within each group, only that the metric of y can be transformed so that this is true for the resulting transformed scores. Without loss of generality, we assume that f is defined so that y^* is standardized in the population, that is $E[y^*] = 0$ and $Var(y^*) = 1$. Hence we assume that there is a continuous scale for “academic achievement” (y^*) for which all within group distributions are normal. Note that neither y nor y^* is assumed to be normally distributed in the combined population of all groups. We elaborate on the conceptual distinctions between these two metrics in Section 1.5 below.

We are interested in the case where neither y nor y^* is observed. Instead, we observe a “coarsened” version of y . This coarsened version, denoted $s \in \{1, \dots, K\}$, is determined by $K - 1$ distinct ordered threshold values, c_1, \dots, c_{K-1} , where $c_{k-1} < c_k$ for all k :

$$s \equiv k \text{ iff } c_{k-1} < y \leq c_k, \quad (2a)$$

where we define $c_0 \equiv -\infty$ and $c_K \equiv +\infty$. Because f is a monotonic increasing function, s is also a

coarsened version of y^* :

$$s \equiv k \text{ iff } c_{k-1}^* < y^* \leq c_k^*, \quad (2b)$$

where $c_k^* = f(c_k)$. Under our assumption of respective normality, the model-implied proportion of observations in category k for group g is therefore

$$\pi_{gk} = \Phi\left(\frac{\mu_g^* - c_{k-1}^*}{\sigma_g^*}\right) - \Phi\left(\frac{\mu_g^* - c_k^*}{\sigma_g^*}\right) = \Pr(c_{k-1}^* < y_{ig}^* \leq c_k^*) \equiv \Pr(c_{k-1} < y_{ig} \leq c_k), \quad (3)$$

where $\Phi(\bullet)$ is the standard normal cumulative distribution function. The aim is to estimate μ_g^* and σ_g^* for each group based on the observed frequencies of members of group g in each of the K ordered proficiency categories.

Equation (3) is an instance of the heteroskedastic ordered probit (HETOP) model. Here we think of each student's ordered proficiency category as the result of a draw from an underlying continuous (normal) distribution of test scores within a group. The HETOP model is an extension of the homoskedastic ordered probit model that allows for heteroskedasticity in the variances of the underlying continuous variable across groups. In the remainder of the paper, we refer to the ordered probit model in which all group variances are assumed equal as the homoskedastic ordered probit (HOMOP) model. The ordered probit model is sometimes referred to as an ordered choice model (Williams, 2009) or as a location-scale model (Cox, 1995; McCullagh, 1980). Most broadly, it is an instance of a generalized linear model that parameterizes the multinomial distribution of observations in each group as cumulative probabilities from a normal density function (Agresti, 2002). Use of a HETOP model allows us to relax the often unrealistic assumption that test scores are homoskedastic across groups and to obtain direct estimates of the within-group standard deviations. To our knowledge, the HETOP model has not been used for the recovery of means and standard deviations from the coarsened data of multiple groups.

Our proposed application and interpretation of the HETOP model is analogous to the ordered

regression model used in the analysis of receiver operating characteristic (ROC) curves, where the model can be interpreted as estimating the mean and standard deviation of unobserved (latent) normal distributions across multiple groups. Tosteson and Begg (1988) demonstrated that the HETOP model generalizes the binormal model for analyzing ROC curves comparing two groups (Dorfman & Alf, 1969) to scenarios with more than two groups. The binormal model has been used previously as a method to estimate the nonparametric V gap statistic between two groups when only coarsened proficiency data are available (Ho & Reardon, 2012; Reardon & Ho, 2015).¹ The HETOP model also generalizes the maximum-likelihood (ML) based estimator of V recommended by Ho and Reardon (2012). It effectively allows for simultaneous estimation of all pairwise V gaps on a common metric for three or more groups.

1.2 HETOP Model Estimation and Identification

Let \mathbf{N} be an observed $G \times K$ matrix with elements n_{gk} containing the counts of observations in group g for which $s = k$; let $\mathbf{P} = [p_1, \dots, p_G]$ be the $1 \times G$ vector of the groups' proportions in the population; and let $\mathbf{n} = [n_1, \dots, n_G]$ be the $1 \times G$ vector of the observed sample sizes in each group, with $N = \sum_g n_g$.² We would like to estimate the vectors $\mathbf{M}^* = [\mu_1^*, \dots, \mu_G^*]^t$, $\mathbf{\Sigma}^* = [\sigma_1^*, \dots, \sigma_G^*]^t$ and $\mathbf{C}^* = [-\infty, c_1^*, \dots, c_{K-1}^*, +\infty]$. In practice, it is preferable to estimate $\mathbf{\Gamma}^* = [\gamma_1^*, \gamma_2^*, \dots, \gamma_G^*]^t$, where $\gamma_g^* = \ln(\sigma_g^*)$. This ensures that the estimates of σ_g^* will all be positive. Following estimation of $\mathbf{\Gamma}^*$, we have $\widehat{\mathbf{\Sigma}}^* = [e^{\widehat{\gamma}_1^*}, \dots, e^{\widehat{\gamma}_G^*}]^t$. Given \mathbf{M}^* , $\mathbf{\Gamma}^*$, and \mathbf{C}^* , and under the assumption of conditional independence of scores within groups, the log likelihood of drawing a sample with observed counts \mathbf{N} is:

¹ The V statistic is a transformation-invariant metric quantifying the nonoverlap between two distributions and is equal to a Cohen's d standardized mean difference when both distributions are normal (Ho, 2009).

² Note that we do not require that $p_g = n_g/N$; that is, the size of the sample in each group need not be proportional to the group's share of the population.

$$\begin{aligned}
L = \ln[P(\mathbf{N}|\mathbf{M}^*, \mathbf{\Gamma}^*, \mathbf{C}^*)] &= \sum_{g=1}^G \left\{ \ln(n_g!) + \sum_{k=1}^K [n_{gk} \ln(\pi_{gk}) - \ln(n_{gk}!)] \right\} \\
&= A + \sum_{g=1}^G \sum_{k=1}^K n_{gk} \ln(\pi_{gk}) \\
&= A + \sum_{g=1}^G \sum_{k=1}^K n_{gk} \ln \left[\Phi \left(\frac{\mu_g^* - c_{k-1}^*}{e^{\gamma_g^*}} \right) - \Phi \left(\frac{\mu_g^* - c_k^*}{e^{\gamma_g^*}} \right) \right],
\end{aligned} \tag{4}$$

where $A = \ln \left(\frac{\prod_{g=1}^G n_g!}{\prod_{g=1}^G \prod_{k=1}^K n_{gk}!} \right)$ is a constant based on the observed counts in \mathbf{N} .

Without constraints on the parameters, the scale of \mathbf{M}^* , $\mathbf{\Gamma}^*$, and \mathbf{C}^* is indeterminate up to a linear transformation. The constraints $\sum_g p_g \hat{\mu}_g^* = 0$ and $\sum_g p_g \hat{\mu}_g^{*2} + \sum_g p_g e^{2\hat{\gamma}_g^*} = 1$ together imply that \mathbf{y}^* has mean 0 and variance 1, as desired. However, these non-linear constraints are not easily implemented in standard software. Instead, it is easier to fit the model subject to two linear constraints on the parameters. As a default we recommend the constraints

$$\begin{aligned}
\mathbf{P}\hat{\mathbf{M}}' &\equiv \sum_{g=1}^G p_g \hat{\mu}_g' = 0 \\
\mathbf{P}\hat{\mathbf{\Gamma}}' &\equiv \sum_{g=1}^G p_g \hat{\gamma}_g' = 0,
\end{aligned} \tag{5}$$

where we use a superscript prime symbol to denote the metric defined by the linear constraints.³

To estimate a homoskedastic ordered probit (HOMOP) model, we impose the additional

³ These specific constraints are not essential; other constraints will identify the parameters, and may be preferable in some settings. The default in many software programs is to define some group r as the “reference group” and to constrain $\mu_r' = 0$ and $\gamma_r' = 0$. These constraints imply that the reference group has a mean of 0 and a standard deviation of 1, with the means and standard deviations of the other groups are then interpreted relative to group r . This is a reasonable default where there is a substantively important reference group and standardization is not needed. It is not the obvious default when there is no substantively important reference group and we would like to estimate each group’s mean and standard deviation relative to the overall population distribution.

constraint that $\hat{\gamma}'_1 = \hat{\gamma}'_2 = \dots = \hat{\gamma}'_G$ before maximizing Equation (4),⁴ so that all groups have identical standard deviations. In some cases, as we discuss below, we may wish to fit a partially heteroskedastic ordered probit (PHOP) model, in which we constrain some subset of the groups to have identical standard deviations, but we allow the others to vary freely. This is achieved by adding to (5) the constraint that the relevant elements of $\hat{\mathbf{\Gamma}}'$ are equal to one another.

We can then maximize Equation (4) subject to the constraints, resulting in estimates $\hat{\mathbf{M}}'$, $\hat{\mathbf{\Gamma}}'$, and $\hat{\mathbf{C}}'$ from which we obtain $\hat{\mathbf{\Sigma}}' = [e^{\hat{\gamma}'_1}, e^{\hat{\gamma}'_2}, \dots, e^{\hat{\gamma}'_G}]^t$. Note that the constraints in (5) (or any set of linear constraints) do not, in general, yield estimates that satisfy the requirement that $\sum_g p_g \hat{\mu}'_g{}^2 + \sum_g p_g e^{2\hat{\gamma}'_g} = 1$. We can, however, standardize the estimates to recover estimates of \mathbf{M}^* , $\mathbf{\Sigma}^*$, and \mathbf{C}^* , using:

$$\begin{aligned}\hat{\mathbf{M}}^* &= \frac{1}{\hat{\sigma}'} \hat{\mathbf{M}}' \\ \hat{\mathbf{\Sigma}}^* &= \frac{1}{\hat{\sigma}'^2} \hat{\mathbf{\Sigma}}' \\ \hat{\mathbf{C}}^* &= \frac{1}{\hat{\sigma}'} \hat{\mathbf{C}}'\end{aligned}\tag{6}$$

where $\hat{\sigma}'$ is an estimate of the population standard deviation in the metric defined by the constraints (the “prime” metric). We show in Appendix A that $\hat{\sigma}'$ can be estimated as:

$$\hat{\sigma}' = \sqrt{\hat{\sigma}_B^2 + \hat{\sigma}_W^2}\tag{7}$$

where

$$\hat{\sigma}_W^2 = \frac{\mathbf{P}\hat{\mathbf{\Sigma}}'^2}{1 + 2\widehat{\omega}_g^2}\tag{8}$$

and

⁴ If we are using the default constraint of $\mathbf{P}\hat{\mathbf{\Gamma}}' = \mathbf{0}$, then this together with the additional homoskedasticity constraint implies the single combined constraint $\hat{\gamma}'_1 = \hat{\gamma}'_2 = \dots = \hat{\gamma}'_G = 0$.

$$\hat{\sigma}_B'^2 = \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \frac{[\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P})]\hat{\boldsymbol{\Sigma}}'^{\circ 2}}{1 + 2\widehat{\omega}_g^2}. \quad (9)$$

In these equations $\widehat{\omega}_g^2$ is the estimated average sampling variance of the $\hat{\gamma}_g'$'s and the “ $(\mathbf{A})^{\circ b}$ ” notation indicates the matrix whose elements are the corresponding elements of matrix \mathbf{A} raised to the power b .

Appendix A shows that for the HETOP model we can use the approximation $\widehat{\omega}_g^2 \approx (2\tilde{n})^{-1}$, where \tilde{n} is the harmonic mean of $n_g - 1$: $\tilde{n} = \left(\frac{1}{G} \sum_g \frac{1}{n_g - 1}\right)^{-1}$. For the HOMOP and PHOP models, $\widehat{\omega}_g^2$ is approximated slightly differently (see Appendix A).

As we noted above, the model can be estimated with different constraints than those in Equation (5), as long as two independent constraints are used. However, if an alternate set of constraints are used it is necessary to transform the resulting estimates to the metric defined by Equation (5) before standardizing using the procedure in Equation (6); we describe the necessary transformation in Online Appendix A. Absent problems maximizing the likelihood function, such as those discussed in Section 1.6, these transformation and standardization procedures will yield the same estimates of \mathbf{M}^* and $\boldsymbol{\Sigma}^*$ regardless of the linear constraints imposed to identify the model.

Maximum likelihood estimation of the HETOP, HOMOP, and PHOP models can be conducted in a number of widely available statistical packages; see Greene and Hensher (2010, p. 179) for a fairly recent list. For all simulations and analyses described in this paper, we carry out the maximum likelihood estimation of (5) using a modification of the `oglm` (“ordinal generalized linear models”) routine (Williams, 2010) written for *Stata* (StataCorp, 2013).⁵

⁵ The modified program is contained in a Stata ado file freely available to readers from the authors upon request.

1.3 Additional Estimands of Interest

Once we have obtained $\hat{\mathbf{M}}^*$ and $\hat{\mathbf{\Sigma}}^*$, estimation of summary statistics like between-group gaps and ICCs is straightforward. First, the achievement gap between any two groups g and h can be computed as the standardized mean difference in y^* between the groups:

$$D_{gh} = \frac{\hat{\mu}_g^* - \hat{\mu}_h^*}{\sqrt{\frac{1}{2}(\hat{\sigma}_g^{*2} + \hat{\sigma}_h^{*2})}} \quad (10)$$

Note that, under the assumption of respective normality, D_{gh} is equal to V , the gap statistic invariant to monotonic scale transformations (Ho & Reardon, 2012).

Second, the ICC (the between-group share of total test score variance) is simply one minus the estimated within-group variance of y^* , because the total variance of y^* is 1:

$$\widehat{\text{ICC}} = 1 - \hat{\sigma}_W^{*2} = 1 - \left[\frac{\mathbf{P}(\hat{\mathbf{\Sigma}}^*)^{\circ 2}}{1 + 2\widehat{\omega}_g^2} \right]. \quad (11)$$

1.4 Computation of Standard Errors

Once we have standardized the estimated group means and standard deviations using Equation (6), we can also compute their standard errors. Because the elements of $\hat{\mathbf{M}}^*$ and $\hat{\mathbf{\Sigma}}^*$ are the products of error-prone estimates of σ' and error-prone elements of \mathbf{M}' and $\mathbf{\Sigma}'$, the standard errors of the elements of $\hat{\mathbf{M}}^*$ and $\hat{\mathbf{\Sigma}}^*$ will depend on the variances and covariances of $\hat{\sigma}'$ and the elements of $\hat{\mathbf{M}}'$ and $\hat{\mathbf{\Sigma}}'$. In Appendix B we derive formulas to estimate \mathbf{V}^* and \mathbf{W}^* , the sampling variance-covariance matrices of $\hat{\mathbf{M}}^*$ and $\hat{\mathbf{\Sigma}}^*$, respectively, when the model is fit with the constraints $\mathbf{P}\hat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\hat{\mathbf{\Sigma}}' = \mathbf{0}$. These derivations take into account the sampling error in $\hat{\sigma}'$.

The standard errors of the gaps described in Equation (10) can be computed from $\hat{\mathbf{M}}^*$, $\hat{\mathbf{\Sigma}}^*$, $\hat{\mathbf{V}}^*$, and $\hat{\mathbf{W}}^*$, as described in Appendix D. The formulas there are generalizations of the formulas used in

Reardon and Ho (2015).

The standard error of the ICC is relatively straightforward to compute once we have $\widehat{\mathbf{W}}^*$. Given Equation (11), the variance of the ICC estimator will be

$$Var(\widehat{ICC}) \approx \left(\frac{1}{1 + 2\widehat{\omega}_g^2} \right)^2 Var \left[\mathbf{P}(\widehat{\boldsymbol{\Sigma}}^*)^{\circ 2} \right] \approx 4 \left(\frac{1}{1 + 2\widehat{\omega}_g^2} \right)^2 \mathbf{P}[\text{diag}(\boldsymbol{\Sigma}^*)] \mathbf{W}^* [\text{diag}(\boldsymbol{\Sigma}^*)] \mathbf{P}^t. \quad (12)$$

Substituting $\widehat{\boldsymbol{\Sigma}}^*$ and $\widehat{\mathbf{W}}^*$ and the appropriate approximation of $\widehat{\omega}_g^2$ (see Appendix A) into Equation (12) yields an estimate of the variance of the ICC estimator.

1.5 A Note on Interpreting the Estimated Parameters

There are two different test score scales relevant to the interpretation of estimated parameters. The first is the continuous scale in which test scores are constructed (i.e., the scale score metric of a test as constructed by a test developer). We denote the scores measured in this metric (the original test metric) with the variable y and denote estimates based on these scores as $\hat{\mu}_k$ and $\hat{\sigma}_k$. The second is the scale of the standardized estimates produced by the HETOP model. We denote the scores measured in this metric with the variable y^* and denote estimates in this metric with a superscript “star” (e.g., $\hat{\mu}_k^*$ and $\hat{\sigma}_k^*$). The estimates in this metric are interpreted relative to a population mean and standard deviation of 0 and 1 respectively. Scores in the prime metric described in Section 1.2 are simply a linear transformation of y^* used in the process of model estimation and are not relevant to the final interpretation of y^* .

If the function f that transforms y into y^* is non-linear, then the group means and standard deviations in the y^* metric will not be linearly related to those in the original y metric. In other words, the target parameters of our application of the HETOP model are not the test score means and standard deviations in the (potentially observed) test score metric of y . Rather, they are the means and standard

deviations in the continuous metric of y^* —the metric in which each group’s distribution is normal and in which the population distribution has mean 0 and standard deviation 1. The key assumption of the model is that such a metric exists. That is equivalent to saying that the group distributions of y (if y could be observed) are “respectively normal,” as defined above.

In some cases, these parameters may be unsatisfying. If, for example, we want to recover means and standard deviations in the reported metric of y (e.g., if we want to recover group-specific mean SAT scores, expressed in the SAT score scale metric), we could do so if two conditions are met. First, we would need to know the threshold scores (in the original metric) used to coarsen the data (that is, we would need to know c_1, \dots, c_{K-1} , where $K \geq 3$). Second, we would have to assume that the group-specific distributions of scores are normally distributed in the original metric (rather than assuming only that they could be normalized by some common transformation f). If these two conditions are met, we could estimate the HETOP model using the frequency counts within each ordered category, as above, except that we would constrain the vector \mathbf{C} to have values equal to the known threshold scores (rather than imposing constraints on the vectors of estimated means and standard deviations). The vectors $\hat{\mathbf{M}}'$ and $\hat{\mathbf{\Sigma}}'$ would then be freely estimated and would be interpretable as group means and standard deviations in the original score metric y . From a practical standpoint, if the group distributions in the original metric are already normal (or nearly normal), the y and y^* estimates will differ only by a linear transformation (or a nearly linear transformation). While the scale of the means and standard deviations will thus differ, auxiliary statistics such as standardized mean differences or the estimated ICC will be unchanged (or nearly unchanged).

When it is reasonable to assume distributions of the original scores y are normal within each group and it is desirable to obtain estimates in that metric, then constraining the thresholds to their known values may be preferable. Unlike physical properties like height or weight, however, it is not clear that there is any natural cardinal scale for cognitive or academic skill (Lord, 1980) or that the test design

principles necessary to support cardinality have been addressed for many common test score scales (Briggs, 2013). In many cases, then, the fixed intervals between established cutscores defined in the original scale score metric might have little theoretical justification or relevance, and may be unnecessarily restrictive. The y^* metric provides a unique metric interpretable in standard deviation units for comparing test score distributions across groups. The y^* metric also preserves the ordinal structure of the observed data, while remaining invariant to plausible monotonic transformations of the original test score scale (i.e., it does not rely on the cardinality of the original test score scale). We therefore use the parameters on the y^* scale as targets for simulation and interpretation.

1.6 Estimation Issues and the Partially Heteroskedastic Ordered Probit Model

The HETOP model will be unidentified if there are groups in which all observations fall in the highest or lowest proficiency category. In this case, the ML estimates of these groups' means will be $\pm\infty$. In such cases, the HETOP model does not have enough information to provide estimates. In other cases, heteroskedastic multinomial or ordered probit models may suffer from fragile identification (Freeman, Keele, Park, Salzman, & Weickert, 2015; Keane, 1992), meaning that although the model is formally identified, the likelihood function may be nearly flat over a range of the parameter space. This may result in a near-singular Hessian, failure of the estimation algorithm to converge, or convergence with very large standard errors. In addition, ML algorithms can sometimes indicate convergence even when the multinomial or ordinal probit model is formally unidentified, due to approximation errors in estimating the likelihood function (Horowitz, Sparmann, & Daganzo, 1982; Keane, 1992).

In our simulations and in applying the HETOP model to real data, we found evidence of such fragile identification in some cases. This occurred when one or more groups had sparse data—for example, when the coarsened data showed all members of a group scoring in the same one or two ordered categories. This condition is unlikely to occur unless more than one of the following conditions

hold: the group has a relatively high or low mean, a small standard deviation, a small sample size, and/or the cutscores are narrowly or unevenly spaced. In such cases, the HETOP algorithm sometimes either failed to converge or converged and returned estimates with very large standard errors for particular groups' parameter estimates (often many orders of magnitude larger than those for other, well-identified groups). In some cases, the algorithm would converge using one set of constraints but not another, or would converge with two different sets of constraints but result in differing estimates of μ_g^* and σ_g^* , suggesting these parameters were at best tenuously identified and not to be trusted.⁶

In such cases, one can drop sparsely populated groups from the model and fit the HETOP model only with groups with sufficient data to identify their parameters. One disadvantage of this is that the standardization procedure we describe will no longer include the full population of interest. An alternate solution is to fit a PHOP (or HOMOP) model instead of the HETOP model, imposing some constraints on the standard deviations of the groups with sparse data. For example, constraining all groups with small sample sizes, or with similar values of some covariate, to have the same standard deviation allows the model to use information from multiple groups to estimate a common standard deviation for those groups. As long as at least some of the constrained groups have sufficient data to identify the parameters, the fragile identification problem may be avoided. We describe simulation analyses of such a model in Section 2.3 below, where we find that the PHOP model often yields a smaller root mean squared error (RMSE) than the HETOP model, even when the groups' true standard deviations are not identical.

⁶ Even in cases where all groups have sufficient data to identify the model parameters, small sample sizes may slow or impede convergence of the ML algorithm, because the likelihood function may be very flat over a wide range of the parameter space. In such cases, we have found that replacing the constraints $\mathbf{P}\hat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\hat{\mathbf{\Gamma}}' = \mathbf{0}$ with a reference group constraint (i.e., constrain $\hat{\mu}_r' = 0$ and $\hat{\gamma}_r' = 0$ where r indicates a reference group) sometimes improved the speed of convergence. In such cases, convergence is improved when the reference group is one with a large sample size and a distribution of frequency counts that is similar to the population distribution. The speed of convergence can also be improved by providing the algorithm with feasible starting values, which can be obtained by using the two-group methods described in Ho & Reardon (2012) to separately estimate each group's mean and standard deviation relative to that of the selected reference group.

2 Evaluating the Performance of the HETOP and PHOP Models Using Simulated Data

We conducted a Monte Carlo simulation to evaluate the accuracy of our proposed use of the HETOP model (and our described standardization procedure) when the data generating procedure matches the model. The first simulation study uses a range of conditions selected to represent those likely to be encountered when analyzing coarsened proficiency data in practice. It builds upon prior simulation studies of HETOP models that sought to recover individual-level parameters rather than group parameters (e.g., Keele & Park, 2006). We focus directly on recovery of the means, standard deviations, and ICCs of the continuous y^* variable after applying our proposed standardization procedure, including evaluation of bias, sampling variability and confidence interval coverage of the estimated standard errors. We also evaluate the performance of the partially heteroskedastic (PHOP) model as a potential way to overcome estimation problems caused by small sample sizes.

2.1 HETOP Model Simulation Conditions and Procedure

We simulated data from populations that differ in the degree to which the true means and standard deviations of test scores vary among groups. We characterize the variation in group means using the ICC (the proportion of total variance in test scores that lies between groups) and the variation in group standard deviations using the coefficient of variation (CV) of group variances (defined as $CV = SD(\sigma^2)/E[\sigma^2]$). We first created four populations, each defined by an ICC (0.05 or 0.20) and a CV (0.0 or 0.3) and containing 100 groups. In each population, each of the 100 groups has one of 10 uniformly-spaced true means and one of 10 uniformly-spaced true standard deviations (when $CV=0$, all groups have identical standard deviations), with the set of means and standard deviations defined so that the population has the desired ICC and CV, and an overall mean of 0 and total variance of 1. We selected the ICC and CV values to correspond roughly to the high and low ends of values reported in prior literature on test score variation (Hedberg & Hedges, 2014; Hedges & Hedberg, 2007) and the real test score data we

analyze later in this paper.

In each of the four populations, we conducted 4 sets of simulations, each defined by groups of a different sample size ($n = 25, 50, 100, 400$). For each of the 16 resulting simulation scenarios defined by the ICC, CV, and group sample size, we generated random samples of size n from each of the 100 groups. Each group's sample was drawn from normal population distributions with means and standard deviations defined by the parameters for each of the 100 groups. We then coarsened the scores four different ways, each time using a different set of cut score locations (set at the 20th/50th/80th; 5th/50th/95th; 5th/30th/55th; and 5th/ 25th/ 50th/ 75th/ 95th percentiles of the population test score distribution, and described as "mid," "wide," "skewed," and "many" cutscores, respectively). The cutscore locations were chosen to be representative of a wide range of conditions found in empirical coarsened test score data (Reardon & Ho, 2015). Finally, we fit both the HETOP and HOMOP model to the coarsened data, and followed the procedures described above to obtain $\hat{\mathbf{M}}^*$, $\hat{\mathbf{\Sigma}}^*$, the estimated ICC, and their standard errors. For each of the 64 scenarios, we repeated this process 1000 times.

Although our primary goal is to assess the performance of the HETOP model, we also fit the HOMOP model to each simulated data set in order to compare the relative performance of the two models. Fitting the HOMOP and HETOP model to data generated from a population that is homoskedastic (CV=0.0) allows us to assess whether using the HETOP model when it is not needed leads to bias or inefficient estimation relative to the more appropriate HOMOP model. Likewise, fitting both models to data generated from a population that is heteroskedastic (CV=0.3) allows us to assess whether and how much the use of the HETOP model improves estimation relative to the HOMOP model.

We evaluated the performance of the HETOP and HOMOP models by computing the bias and RMSE of the estimated means, standard deviations, and ICCs. For the means and standard deviations we focused primarily on the aggregate bias and RMSE (averaged across all $G = 100$ groups) for each estimator $\hat{\theta}$ (where θ could be a mean or a standard deviation) by computing:

$$Bias_{\hat{\theta}} = \frac{1}{R * G} \sum_{r=1}^R \sum_{g=1}^G (\hat{\theta}_{gr} - \theta_g) \quad (13)$$

$$RMSE_{\hat{\theta}} = \sqrt{\frac{1}{R * G} \sum_{r=1}^R \sum_{g=1}^G (\hat{\theta}_{gr} - \theta_g)^2}$$

where R is the number of converged replications (usually 1000), $\hat{\theta}_{gr}$ is the estimate for group g in replication r and θ_g is the true value. For the ICC estimates, \widehat{ICC} , bias and RMSE were computed as:

$$Bias_{\widehat{ICC}} = \frac{1}{R} \sum_{r=1}^R (\widehat{ICC}_r - ICC) \quad (14)$$

$$RMSE_{\widehat{ICC}} = \sqrt{\frac{1}{R} \sum_{r=1}^R (\widehat{ICC}_r - ICC)^2}.$$

To evaluate the accuracy of our formulas for the standard errors (SE) of group means and standard deviations we computed the average ratio of the median⁷ estimated SE of a parameter to its empirical SE (the standard deviation of the sampling distribution of the parameter) across all $G = 100$ groups in a condition:

$$SE\ Ratio_{\hat{\theta}} = \frac{1}{G} \sum_{g=1}^G \frac{Median(\widehat{SE}_{\hat{\theta}_{gr}})}{SD(\hat{\theta}_{gr})}. \quad (15)$$

To evaluate the accuracy of the standard error formula of the estimated ICC's, we compute the ratio of the median estimated SE of the ICC to its empirical SE:

$$SE\ Ratio_{\widehat{ICC}} = \frac{Median(\widehat{SE}_{\widehat{ICC}_r})}{SD(\widehat{ICC}_r)} \quad (16)$$

⁷ We used the median rather than the mean estimated standard error to reduce the impact of extreme standard error estimates, primarily in conditions with small sample sizes and wide or skewed cutscores.

where $SD(\hat{\theta}_{gr})$ and $SD(\widehat{ICC}_r)$ are the observed standard deviations of the sampling distributions of the relevant parameter estimates across the R replications for a given condition. If our SE formulas in Appendix B are accurate, we expect the ratios in (15) and (16) to be close to 1. We also computed the 95% confidence interval (CI) coverage rates for each parameter, computed as the proportion of cases for which $|\hat{\theta}_{gr} - \theta_g| < 1.96 * \widehat{SE}_{\theta_{gr}}$ or $|\widehat{ICC}_r - ICC| < 1.96 * \widehat{SE}_{ICC_r}$. If the estimates are biased, the CI coverage rates will not equal 95%, however, even if the standard error formulas accurately reflect sampling variability.

Finally, we present results describing the loss of efficiency (in terms of increased sampling variance) when estimating group means and standard deviations from coarsened rather than full data. For each condition we estimate relative efficiency as the average efficiency ratio across groups. We define this as the average (across groups) of the ratio of the observed sampling variance of the target parameter in the simulations (using coarsened data) to its theoretical sampling variance if it were estimated from continuous data:

$$Average\ Efficiency\ Ratio_{\theta} = \frac{1}{G} \sum_{g=1}^G \frac{\widehat{Var}(\hat{\theta}_{gr})}{\tau_{\theta_g}^2}, \quad (17)$$

where θ is either a mean or a standard deviation, $\widehat{Var}(\hat{\theta}_{gr})$ is the observed variance in estimates of the target parameter across the R replications and $\tau_{\theta_g}^2$ is the theoretical sampling variance of the estimator based on continuous data (when θ is the mean, μ_g^* , then, $\tau_{\theta_g}^2 = \sigma_g^2/n_g$; when θ is the standard deviation, σ_g^* , then, $\tau_{\theta_g}^2 = \sigma_g^2 * (2[n_g - 1])^{-1}$). An efficiency ratio of 1.0 would indicate that estimates based on coarsened data have the same sampling variance as estimates based on continuous data; efficiency ratios larger than 1.0 indicate there is greater variability in estimates based on coarsened data. The efficiency ratio can be interpreted as the ratio by which the sample size would need to be increased to estimate the parameters from coarsened data with the same precision as if the parameters were estimated from

continuous data.

Because the 100 true group means and standard deviations were held constant across the 1000 replications within a given scenario, we also examine the bias, RMSE and SE performance for individual groups within a particular condition when relevant. Online Appendix C contains detailed tables of all aggregate bias, RMSE and SE results. We used the constraints $\mathbf{P}\hat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\hat{\mathbf{\Gamma}}' = \mathbf{0}$ to identify the model, and the ML algorithm converged in all but 5 of the total 128,000 replications.

2.2 HETOP Model Simulation Results

2.2.1 *Recovery of Means*

The aggregate bias of means estimated with the HETOP model was indistinguishable from 0 for all conditions, and the bias for individual groups was also indistinguishable from 0 in almost all of the scenarios we explored. The one exception was in the small sample ($n=25$) simulations with a large ICC (0.20), large CV (0.3) and skewed or wide cutscores; in these cases, we detected non-zero bias for some groups, though the bias was very small, nearly always less than 0.05 standard deviation units for any given group. Moreover, not only was this bias small in absolute terms, but it was also very small in relation to the aggregate RMSE of the estimated means (which in this case was on average approximately ten times larger than the largest bias we observed). We do not show these results for parsimony. The precision of the estimated means varied primarily as a function of sample size, although when sample sizes were small the sampling variance was modestly affected by the location of the cutscores; sampling variance was consistently lowest for estimates based on the “many” cut score condition, as one would expect given the light degree of coarsening.

2.2.2 *Recovery of Standard Deviations*

The top panel of Figure 1 shows the average bias in standard deviation estimates across all groups and all replications for each condition. This figure illustrates that there is some negative bias in the

standard deviation estimates from the HETOP model and that the bias is primarily a function of sample size that is exacerbated when cutscores are skewed or wide. Note the average bias for standard deviations is quite small compared to the true standard deviation of scores, typically less than 1% of the size of the true standard deviations, except when sample sizes are less than 50 and the cutscores are skewed or wide (the average standard deviation is approximately 0.89 when ICC = 0.20 and 0.97 when ICC = 0.05, and the largest absolute bias in any condition is approximately 0.045).

[Figure 1 here]

When CV = 0 and the HOMOP model is the correct model, the top panel of Figure 1 indicates a very slight negative bias in HOMOP standard deviation estimates that generally approaches 0 with increasing sample size more quickly than the corresponding HETOP estimates, particularly with skewed or wide cutscores. That is, when the group distributions are truly homoskedastic and the coarsening is done sub-optimally, the HOMOP model produces less biased estimates of standard deviations than the HETOP model. Note, however, that the HOMOP model produces modest positive bias in the standard deviation estimates in the CV = 0.3 conditions (where the HOMOP model is not the correct model), particularly when sample sizes are large. Given the misspecification of the model, such bias is not surprising.

The bottom panel of Figure 1 depicts the RMSE of standard deviation estimates across conditions as defined in Equation (13). The results in Figure 1 together suggest that when the data are homoskedastic, the HOMOP model is generally (and unsurprisingly) preferable to the HETOP model. If data are heteroskedastic, however, the HOMOP model will systematically over/underestimate individual group standard deviations, with bias inversely related to the true standard deviation. Nonetheless, if one wishes to minimize RMSE, it may still be better to use a HOMOP model if sample sizes are small. In the scenarios shown in Figure 1 the HOMOP model generally yields a smaller RMSE than the HETOP model for scenarios with $n < 100$. The sample size at which the HOMOP model is preferable to the HETOP model (in terms of RMSE) will be a function of a number of factors, particularly the CV of group variances and

the location of the cutscores. We investigate this bias/variance tradeoff further in Section 2.3.

Figure 2 provides more detail on the systematic patterns of bias in group standard deviation estimates from the HETOP model, showing the bias in standard deviation estimates as a function of the true population means and standard deviations for the condition in which $ICC=0.20$ and $CV=0.3$. Each panel in Figure 2 shows the bias as a function of groups' true mean and standard deviation for a given sample size and cut score condition, with the x-axis indicating group means and y-axis indicating group standard deviations. The figure makes clear that the bias in estimated standard deviations varies with a group's true mean and standard deviation, when cutscores are skewed or wide and sample sizes are small. The top left panel, for example, shows that nearly all group standard deviation estimates are negatively biased when $n=25$ and cutscores are skewed, but the bias is largest for groups with larger true means and smaller true standard deviations. This pattern is a result of the loss of information due to the coarsening of the data. If a group's true standard deviation is small and its true mean is high, then when the cutscores are skewed the coarsening leads to observed data with little information (most cases will fall in the top category, providing little information about the group's standard deviation) and to underestimation of standard deviations. In larger samples, coarsening leads to much less consequential loss of information, however, as is evident in the bottom ($n=400$) row of panels, where absolute bias is less than 0.01 for all but one group across all cutscore conditions. Although not pictured, the pattern of small negative bias in small groups is similar when $CV=0.0$.

[Figure 2 here]

2.2.3 *Recovery of ICC*

Figure 3 shows the bias in ICC estimates across conditions. The ICC estimates are upwardly biased, particularly when sample sizes are small; this is partly a result of the small negative bias in standard deviations in these cases (see Figure 1). The bias in the HETOP ICC estimates does not appear to depend on the true CV, but is modestly larger when the true ICC is larger. When $CV=0$, the HETOP and

HOMOP estimates are similarly biased in most cases, though the bias is slightly lower in the HOMOP model when the cutscores are skewed or wide. When the data are heteroskedastic ($CV = 0.3$), the HOMOP ICC estimates are biased even when $n = 400$, due to the misspecification of the model. In all cases, the bias is largest when cutscores are skewed or wide. Overall, however, the bias in ICCs is relatively small, generally less than 0.01 unless $n = 25$ or the cutscores are skewed.

[Figure 3 here]

2.2.4 Accuracy of Standard Errors

The accuracy of the standard errors in the simulations was similar across ICC and CV conditions. For parsimony and to limit sampling variability, Table 1 shows the SE ratios and CI coverage rates averaged across the 4 combinations of ICC (.05 and .20) and CV (0.0 and 0.3) conditions. Table 1 indicates that estimated SEs and CIs for all three parameters were accurate with moderate and large sample sizes ($n = 100$ or more), but less accurate with smaller sample sizes. SEs and CIs were least accurate with small sample sizes when cutscores were widely spaced. In such cases, the approximations used to derive the standard error formulas (Appendix B) appear to break down.

[Table 1 here]

2.2.5 Efficiency of Estimators

Figure 4 presents the average efficiency ratio across all 100 groups for each condition when using the HETOP model. The top panel shows average efficiency ratios for estimated means while the bottom panel shows the average efficiency ratios for the standard deviations; each panel represents a different ICC and CV condition while each line represents a different cutscore condition. For the means, the loss of efficiency is moderate and depends primarily on the cutscore locations, with the greatest loss of efficiency when the cutscores are skewed low. Within any combination of cutscores, CV, and ICC, the relative loss of precision is largest when samples are small. The average efficiency ratio for estimated means across all groups and conditions is 1.36, ranging from a minimum of 1.06 (in the case where there

are many cutscores, a CV of 0.3, an ICC of 0.05, and $n = 400$) to a maximum of 2.49 (in the case where the cutscores are skewed low, the CV is 0.3, the ICC is 0.20, and $n = 25$). This indicates that in some conditions, the coarsening of the data results in very little loss of precision, while in others (very small group sizes and skewed cutscores) the loss of precision is more substantial.

[Figure 4 here]

The efficiency loss with respect to estimating group standard deviations is larger than when estimating means, but again, the efficiency ratio varies considerably depending on cutscore locations and sample size. The skewed low cutscore condition is consistently the least efficient, and the many cutscores condition is consistently the most efficient, with average efficiency ratios of 3.53 and 1.37, respectively, averaging across all CV's, ICC's, and group sizes. The efficiency ratio in the wide cutscore condition appears to be most dependent upon group sample sizes: when group sample sizes are 25, the average efficiency ratios range from 3.12 to 3.93 across the ICC and CV conditions; they are half as large (1.60 to 1.82) when sample sizes are 400.

2.3 PHOP Model Simulation Conditions and Procedure

When the data are truly homoskedastic, the simulation results in Section 2.2 show that a HOMOP model with all group standard deviations constrained to equality performs better than a fully heteroskedastic model. However, the results also suggest that in some truly heteroskedastic cases with small sample sizes, the HOMOP model may be preferable, as reductions in RMSE could outweigh increases in bias for group standard deviation estimates. In the simulations above, however, all groups in a given simulation scenario had the same sample sizes, a condition that may not often hold in practice.

Anticipating contexts in which sample sizes across groups differ, we evaluate the performance of a partially constrained heteroskedastic ordered probit model (PHOP) in which standard deviation estimates for small groups are constrained to equality while those for large groups are freely estimated.

Because the efficiency/bias tradeoff implicit in constraining group standard deviations will depend on how much true variation in standard deviations there is, we conduct these simulations in populations with different degrees of heteroskedasticity (CV).

We follow the same general simulation methodology as outlined in Section 2.1, but with the following modifications. We generate data from five populations, each with ICC = 0.20, and one of five different CVs of group variances (0.0, 0.1, 0.2, 0.3 or 0.4). Each population contains 36 group types whose means and standard deviations are bivariate uniformly distributed with values set to produce the defined ICC and CV. For each group type, we draw 7 random samples of sizes $n=25, 50, 75, 100, 150, 200,$ and 400 from normal distributions defined by each of the 36 group mean/standard deviation values. This yields $7 \times 36 = 252$ groups, one of each combination of mean, standard deviation and sample size. We then coarsen these scores four separate times, using the same cutscores described above (i.e., the “mid”, “many”, “skewed”, and “wide” conditions). For each coarsened sample, we then fit 8 different models: the HETOP and HOMOP models, as well as PHOP models where groups with sample sizes less than or equal to either 25, 50, 75, 100, 150 or 200 were all constrained to be equal. We repeated this process for 1000 replications. We used the constraints $\mathbf{P}\hat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\hat{\mathbf{\Gamma}}' = \mathbf{0}$ to identify the model. The ML algorithm converged in all 40,000 replications using the “many” cuts cores, and failed to converge in nine, 13 and 271 of the 40,000 replications using the “mid”, “skewed” and “wide” cutcores respectively.

2.4 PHOP Model Simulation Results

In general, results for the bias, RMSE and standard errors were similar for the overlapping HETOP and HOMOP conditions here and in Section 2.2, suggesting that the conclusions above remain largely unchanged for conditions with groups of varying sample sizes. However, in some cases there was less bias in HETOP standard deviation estimates for groups with small sample sizes in the simulations with a range of group sizes. This may result from the fact that the overall standard deviation ($\hat{\sigma}'$) is more accurately

estimated when there are some groups with large sample sizes, so that less bias is introduced when we divide by this estimated standard deviation to obtain the $\hat{\sigma}^*$ estimates. For the PHOP models, the accuracy of the estimated standard errors was very good: the ratio of median estimated to empirical standard errors was close to 1 in all cases.

Our primary motivation for testing the PHOP models was to assess whether they reduce the aggregate RMSE of group standard deviation estimates. Estimating a single pooled standard deviation estimate across small groups will yield more precise (but potentially biased) estimates of small groups' group standard deviations; if the increase in bias is outweighed by the reduction in error, the PHOP model may be preferred. Hence our discussion of the results in this section focuses on the RMSE of the standard deviations of groups of various sizes.

Figure 5 displays the RMSE of group standard deviation estimates for different PHOP models in the condition in which $CV = 0.2$ and data were coarsened with the "mid" cutscores. Each panel of the figure displays the aggregate RMSE of standard deviation estimates for a different PHOP model (e.g., PHOP25 is a model in which group standard deviations are constrained to equality for groups with $n \leq 25$); each panel also includes results for the HOMOP (dotted line) and HETOP (dashed line) models, which are the same across panels as they are not affected by the sample size threshold used in the PHOP model. We show RMSE disaggregated by group size here (unlike in Figure 1) because the PHOP model treats groups of different sizes differently by design. Figure 5 shows that the RMSE of constrained group standard deviation estimates in the PHOP model are nearly identical to HOMOP model RMSEs while the unconstrained group standard deviation estimates are nearly identical to the HETOP RMSEs. The pattern of constrained and unconstrained group RMSEs tracking the HOMOP and HETOP model results was consistent across CV values (not shown).

[Figure 5 here]

Figure 5 suggests there will be an optimal sample size threshold at which to constrain standard

deviation estimates to minimize the overall RMSE. Figure 6 displays information useful for determining such a threshold for each CV -by-cutscores condition. Each panel of Figure 6 shows RMSE of group standard deviations (aggregated across all sample sizes) for each model type and each CV condition. The upper left panel, for example, shows the results for the “many” cutscore condition, and includes a line for each CV condition. For a given CV and cutscore condition, an optimal threshold can be identified by finding the model that minimizes RMSE for the corresponding line. When the true CV is 0, the HOMOP model minimizes RMSE in all conditions. When the true CV is 0.2, the optimal models (among those we tried) would be PHOP50, PHOP100, PHOP75, and PHOP150 for the many, mid, wide and skewed cutscore conditions, respectively. Although these results do not cover all possible combinations of ICC, CV, and cutscore locations, they are suggestive about the conditions under which a PHOP model would minimize the RMSE of group standard deviation estimates. In analysis of real data, analysts will know the location of the cutscores, the number of groups, and the group sizes; they may also have information about the range of plausible values of the ICC and CV. These could be used to conduct customized simulations of the type we show here to make an informed decision about the optimal HETOP/PHOP/HOMOP model to select to minimize RMSE, if that is their goal.

[Figure 6 here]

2.5 Summary of Simulation Analyses

These simulations demonstrate that the HETOP model works well when the model matches the data generating process. Unbiased and precise recovery of standard deviations generally requires group sizes of 100 or more. Figure 1 suggests that in some cases where sample sizes are small, a homoskedastic model may produce more efficient (although biased) standard deviation estimates even if the data are truly heteroskedastic. The results in Section 2.4 suggest that using a PHOP model, which constrains small groups to have equal standard deviation estimates, improved the efficiency of standard deviation

estimates with only a relatively small increase in average bias, thus reducing the RMSE. Although the optimal group size at which to constrain the group standard deviations to be equal is not a priori clear in any given scenario, the results above suggest that the analyst may be able to make an informed choice to achieve a roughly optimal model.

3 Application of the HETOP Model to Real Data

The simulations in Section 2 indicate that the HETOP model accurately recovers means, standard deviations, and ICCs from coarsened data across a range of scenarios when sample sizes are moderately large and the group distributions are normal. This section analyzes 18 sets of real test score data to investigate whether means and standard deviations can be recovered from real coarsened test score distributions. To carry out these analyses, we selected datasets for which we had access to both the coarsened proficiency data and the scale scores (the uncoarsened, continuous data) for each student: 10 datasets from a mid-size state's testing program and eight datasets from the state NAEP administrations in 2009 and 2011. In effect, these analyses assess whether the actual test score distributions in these 18 cases satisfy the respective normality assumption of the HETOP model.

3.1 Data

The first eight datasets contained student-level records for the 2009 and 2011 Grade 4 and 8 Main NAEP mathematics and reading administrations, with each dataset containing scores for a single year-by-grade-by subject combination (e.g., 2009 Grade 4 math scores constitute one dataset). The groups in these datasets were states, and the aim was to estimate the means and standard deviations of state test score distributions with the HETOP model. Hence there were 50 groups in each of the NAEP

datasets, with a median group (state) sample size of 3,050 across all eight datasets.⁸

The other 10 datasets consist of mathematics and reading test scores from a medium-sized state for a cross-section of approximately 90,000 students in grades 4 through 8 during the 2005-2006 school year. Each dataset contained student-level scores from a single grade-by-subject combination, with scores grouped at the school level, so that the target estimates of interest were the school means and standard deviations of test scores for given grade levels. Across the 10 state datasets the number of groups (schools) ranged from 428 to 1,244 and the median group (grade within school) sample size ranged from 70 to 194. Both the NAEP and State testing programs use three unique cutscores in each grade and subject level to classify students into one of four ordered proficiency categories. Online Appendix Table D1 provides detailed descriptive information about the 18 datasets.

3.2 Comparison of HETOP and Uncoarsened Estimates

If the test scores in these 18 datasets are respectively normal, and if the group sample sizes are large and the cutscores are well placed, the HETOP model should return precise, unbiased estimates of the group means and standard deviations in the continuous metric of y^* , as our simulation suggests. If, additionally, the function relating the reported scale scores to the metric in which the distributions are normal is linear then group means and standard deviations based on the student-level scale scores should be perfectly correlated (within the limits of sampling variability) with the group means and standard

⁸ NAEP is administered to a sample of students in the nation, and special scoring and scaling techniques result in “plausible values” (e.g., Mislevy, Johnson, & Muraki, 1992) instead of individual scores. For each of the eight year-grade-subject combinations in the dataset, we had five plausible values for each student. To generate a dataset with a single score for each student that could be used to compare with HETOP estimates, we created a synthetic dataset using the first set of plausible values for all students. In order to avoid complications from comparisons using multiple plausible values and sampling weights, we generated an artificial sample for each state using the following procedure. We drew a random sample with replacement from the set of non-missing first plausible values, with probability of selection proportional to original sampling weights. This created a random sample from a population defined by the (weighted) observed first plausible values in each state. We sampled N_g values for each state, where N_g was the original number of unique students with non-missing data in state g . Using NAEP’s actual proficiency cutscores for each subject/grade combination, we calculated the number of students in the synthetic sample scoring in each proficiency category.

deviations based on fitting the HETOP model to the coarsened proficiency data. This suggests we could examine the correlation between HETOP estimates and estimates based on observed scale scores to assess the extent to which the empirical test score distributions satisfy the respective normality assumption of the model.

An imperfect correlation, however, might result not only from a failure of respective normality, but also might arise if a) the function f is not linear; or b) the estimates are imprecise because sample sizes are not large or the cutscores are not sufficiently informative. The first condition will lead to a nonlinear association between the two sets of estimates. The second will produce a noisy association. To assess the respective normality assumption in real test score data, then, we must determine whether the less-than-perfect correlation between the HETOP estimates and the estimates based on the observed scale scores can be explained by the error that comes from coarsening and/or the non-linearity of f . We describe our approach to doing this below. To the extent that these factors do not explain an observed correlation less than one, the test score distributions are not respectively normal.

First, we estimated group means and standard deviations based on the original student-level scale scores, using traditional estimators of means and standard deviations. We refer to these as the “original” scale score estimates. Second, to model the data that researchers may be limited to in practice, we coarsened the scale scores according to the operational NAEP and State cutscores. We then used the HETOP model to estimate means and standard deviations based on the coarsened frequency counts. We refer to these as the “H4” estimates because they are based on four proficiency categories. The correlation between these two sets of estimates will be degraded by imprecision due to the coarsening and by non-linearity in f .

Third, to generate HETOP estimates less affected by loss of information due to coarsening, we coarsened each dataset a second time, using 19 equal-interval cutscores that classified students into 20 “proficiency” categories (instead of four). We then estimated means and standard deviations for each

group with a HETOP model using the 20 observed frequencies for each group. We refer to these as “H20” estimates, because they are based on 20 proficiency categories.⁹

Finally, we estimate a function f^* that, when applied to the observed student-level scale scores simultaneously renders all of the within-group distributions as nearly normal as possible. We estimate f^* from the mapping between the 19 cutscores estimated from the H20 model (i.e., $\hat{c}_1^*, \hat{c}_2^*, \dots, \hat{c}_{19}^*$) and their corresponding values on the reported score scale (c_1, \dots, c_{19}). We estimate a monotonic function that goes through these 19 points (so that $\hat{f}^*(c) = \hat{c}^*$); this function will closely approximate a function that renders the within-group distributions as nearly normal as possible. We then apply this transformation to the observed student-level scale scores, resulting in transformed scores, \hat{y}^* , for each student. If test scores are respectively normal, this transformation should render the group score distributions normal; the group means and standard deviations of \hat{y}^* will be linearly related, within sampling variability, to those estimated from the HETOP model applied to coarsened data. We refer to group means and standard deviations based on these normalized \hat{y}^* scores as “transformed” estimates. The procedure used to estimate f^* is described in Online Appendix E.

We calculated Pearson correlations between these four sets of estimates for each of the 18 datasets. These correlations are summarized in Table 2, which presents the average, minimum and maximum correlation among the estimates for the NAEP and State datasets separately (correlations for each of the 18 data sets are in Online Appendix Table D2). Column (1), for example, summarizes correlations between means estimated based on the H4 and uncoarsened original scale scores, while column (5) summarizes the corresponding correlations between the standard deviation estimates.

[Table 2 here]

As mentioned above, these correlations may be less than 1.0 even if score distributions are

⁹ We could have used more than 20 categories, but given the finite number of possible scale scores and size of the groups, additional categories add vanishingly little additional information.

respectively normal. If the test score data are respectively normal, however, then we expect the correlations in columns (2) and (6) to be near 1.0, because these correlations are based on estimates that adjust for a lack of normality of the reported scale score metric (transformed) and the error due to coarsening into only 4 categories. Indeed, the average correlations between these estimates are uniformly near 1 for all datasets (the lowest correlation across both columns (2) and (6) is 0.979), suggesting both that the data are respectively normal and that the H20 model accurately recovers the group means and standard deviations in the y^* metric.

To evaluate whether the original scale scores are reported in the metric in which they are normal in each group, we examine two sets of results. First, we inspect the correlations between the original and transformed estimate in columns (3) and (7). If the original test score data were reported in the normal metric, we would expect these correlations to be close to 1, because \hat{f}^* would be linear. Second, we plot the function \hat{f}^* in each case to examine its linearity directly. The correlations are near 1 for the means, but lower (as low as 0.83 in one case) for the standard deviations. The plots of \hat{f}^* (in Online Appendix Figure E1) show very slight nonlinearity in most cases. Both of these patterns indicate that while the original test score scales are generally not one in which the distributions are as near to normal as possible, the original scales are not very different from such a scale. The modest departure from normality appears to cause more discrepancy in the estimated standard deviations (average correlations of 0.955 and 0.932 for NAEP and State) than in the estimated means (average correlations of 0.999 for both the NAEP and State scales).

Finally, it is useful to compare columns (2) and (4) and columns (6) and (8); this comparison indicates the extent to which coarsening into 4 rather than 20 categories reduces the precision of the estimated means and standard deviations. The correlations between the H20 and transformed estimates of means (column 2) are generally only modestly larger than those between the H4 and the transformed estimates (column 4). In the case of the estimated standard deviations, however, coarsening substantially

degrades precision: the correlations in column (8) are much lower than in column (6). This is consistent with our simulation results showing that the HETOP model more reliably estimates group means than standard deviations, particularly when group sizes are small and cutscores are not optimally located, as is the case in the state data sets.

These analyses suggest the assumption of respective normality of test score distributions is reasonable in the datasets we examined, which include both school-level and state-level groups. Moreover, the reported scale scores in these datasets appear to be in a metric that is very close to the latent metric in which the means and standard deviations are estimated by the HETOP model. This may not be true for all empirical test score distributions, of course; it would be useful to test in other cases where continuous scores are available.

4 Discussion

This paper introduces a method for estimating the means and standard deviations of continuous test score distributions in multiple groups using only coarsened proficiency data. Through simulations and real data analyses, we demonstrate that accurate estimation of means and standard deviations of test score distributions for multiple groups (states, districts, schools, etc.) is possible under a wide range of scenarios, with modest loss of efficiency, particularly when sample sizes are larger than 50 and when the cutscores are not highly skewed. The analyses also showed that estimates of secondary statistics such as the ICC can be recovered accurately, with slight positive bias when group sizes are small. While estimates of group standard deviations were accurate across all conditions with larger sample sizes, there was evidence of small negative bias in some conditions with smaller sample sizes, particularly when the location of cutscores used to coarsen the data are unequally and/or unevenly spaced, thus providing relatively little information about the original distributions. The bias was very small when group sample sizes were 100, and modestly larger with small samples of size 25, though the average bias was never

sizeable compared to the true standard deviations or the sampling variance of the estimates. Our analyses of real test score datasets suggest the primary assumption of respective normality is reasonable for these particular test scores and likely those developed under similar conditions. Further simulation studies to evaluate the methodology across a wider range of conditions, including those where data are not respectively normal, would be a useful extension to this work.

The simulation results and real data analyses suggest a few common considerations for researchers to attend to when applying the HETOP model in practice. First, because the quality and reliability of HETOP estimates (particularly for group standard deviations) depend primarily on group sample sizes and cut score locations, an inspection of the overall proportion of students within each proficiency category and the proportion of groups with zero observations in one or more categories can be useful indicators of potential problems. Other indicators include models that either will not converge, are slow to converge, or converge but produce abnormally large standard errors. In these cases, our simulations and other work with real test score data suggest a PHOP model is a good way to overcome some data limitations and is generally preferable to a HOMOP model unless the assumption of homoskedasticity is defensible.

In fitting the PHOP model, the analyst must determine a sample size threshold below which to impose the homoskedasticity constraint. This choice can be guided by knowledge of the cutscore locations, the number and size of groups, and prior research that provides information about plausible values of the ICC and CV. When the CV of group variances is approximately 0.2 (roughly the average value observed in the data we analyzed), constraining the standard deviations of groups smaller than 100 is generally near optimal in our simulations. Of course, RMSE need not be the only criterion used to determine the best model. For analysts who are less willing to tolerate bias than error variance, a smaller constraint threshold would be preferable, and vice versa. In addition, if group means are of primary interest, the choice of a PHOP model will matter little; if standard deviations are of interest, the bias-

precision tradeoff is more salient. Further development of practical model fit statistics and diagnostics that can inform PHOP model selection are an important direction for future research.

One benefit of the PHOP model is that it improves estimation for small groups, particularly when cut score locations are sub-optimal. The challenges for estimation posed by small sample sizes or extreme cutscores could also be addressed with alternative estimation strategies or frameworks, such as Bayesian or random-effects models. It is possible to estimate a mixed-effects HETOP model (see for example, Gu, Fiebig, Cripps, & Kohn, 2009; Hedeker, Demirtas, & Mermelstein, 2009) from which one could obtain shrunken estimates of group means and standard deviations. These Bayesian estimates would have smaller RMSE than our ML estimates, but would also contain more bias. The decision of whether to prefer more-biased, lower-RMSE shrunken estimates over less-biased, higher-RMSE ML estimates depends on how one wants to use the resulting estimates. If the estimates will be used as outcome variables in subsequent models or as descriptive statistics, the (less biased) ML estimates may be preferable to the (more biased) shrunken estimates. If the estimates will be used as predictor variables in subsequent models, however, the shrunken estimates may be preferable (although in this case they should, in principle, be shrunken to their mean conditional on the other covariates to be used in the model). Shear, Castellano, and Lockwood (2016) present some preliminary comparisons of these two approaches in the context of coarsened test score data, but additional work exploring the potential benefits of Bayesian HETOP models would be very useful.

In our discussion here we have ignored the potential effects of measurement error. If we think of the continuous scores in the y metric as containing measurement error, then the key assumption of the HETOP model is that the observed, error-prone test score distributions are respectively normal. Given this assumption, estimation proceeds as we describe it above, and the resulting estimates are understood as means and standard deviations of the error-prone scores in the y^* metric. To recover means and standard deviations of true scores in the y^* metric, one would need information about the reliability of

the test scores in that metric. Although this is not identical to the reliability of scores in the metric y (the metric reported by test score developers), unless the function f is linear, Reardon and Ho (2015) show that using published reliabilities to adjust group means and standard deviations on a transformed scale generally produces only trivial bias, given that widely-used standardized tests typically have high reliability. When reliability is high, distortions of measurement error due to the transformation function f are trivial unless f is extremely non-linear. As a result, standard measurement error adjustments, based on published reliabilities of scores in the y metric, can be made to yield estimates of groups' true test score means and standard deviations in the y^* metric.

Finally, as mentioned above, these methods are applicable whenever data can be conceptualized as coarsened: the result of some form of polychotomization, censoring, binning, or rounding. In the case of aggregate proficiency data, such as that contained in the *EDFacts* database, such a model is clearly applicable, and our results show that the HETOP model can provide estimates of means and standard deviations that can overcome some of the limitations with such data as described by Ho (2008) and others. In the case of AP exams, where scores are only reported on a 1-5 ordinal scale, one might still presume the existence of a continuous underlying variable of which the observed scores are a coarsened version. In such cases, our methods provide a way to estimate the distributions of this underlying continuous variable in multiple groups. Ordinal data of many kinds—from Likert scale survey data to Apgar scores and from discrete levels of educational attainment to demographic age or income bins—can be thought of as representing coarsened versions of latent continuous variables. In many of these cases, the methods described here could be usefully applied to estimate moments of group distributions.

References

- Agresti, A. (2002). *Categorical data analysis* (2nd ed.). Hoboken, N.J.: John Wiley & Sons.
- Alvarez, R. M., & Brehm, J. (1995). American ambivalence towards abortion policy: Development of a heteroskedastic probit model of competing values. *American Journal of Political Science*, *39*(4), 1055–1082.
- Briggs, D. C. (2013). Measuring growth with vertical scales. *Journal of Educational Measurement*, *50*(2), 204–226.
- Casella, G., & Berger, R. L. (2002). *Statistical inference* (2nd ed.). United States: Duxbury.
- Cizek, G. J. (Ed.). (2012). *Setting performance standards: Foundations, methods, and innovations* (2nd ed.). New York, NY: Routledge.
- Cox, C. (1995). Location—scale cumulative odds models for ordinal data: A generalized non-linear model approach. *Statistics in Medicine*, *14*(11), 1191–1203. <http://doi.org/10.1002/sim.4780141105>
- Dorfman, D. D., & Alf, E. (1969). Maximum-likelihood estimation of parameters of signal-detection theory and determination of confidence intervals—rating-method data. *Journal of Mathematical Psychology*, *6*(3), 487–496.
- Freeman, E., Keele, L., Park, D., Salzman, J., & Weickert, B. (2015, August 14). *The plateau problem in the heteroskedastic probit model*. Retrieved from <http://arxiv.org/abs/1508.03262v1>
- Goodman, L. A. (1960). On the exact variance of products. *Journal of the American Statistical Association*, *55*(292), 708–713. <http://doi.org/10.1080/01621459.1960.10483369>
- Greene, W. H., & Hensher, D. A. (2010). *Modeling ordered choices: A primer*. New York: Cambridge University Press.
- Gu, Y., Fiebig, D. G., Cripps, E., & Kohn, R. (2009). Bayesian estimation of a random effects heteroscedastic probit model. *Econometrics Journal*, *12*(2), 324–339. <http://doi.org/10.1111/j.1368-423X.2009.00283.x>
- Hedberg, E. C., & Hedges, L. V. (2014). Reference values of within-district intraclass correlations of academic achievement by district characteristics: Results from a meta-analysis of district-specific values. *Evaluation Review*, *38*(6), 546–582. <http://doi.org/10.1177/0193841X14554212>
- Hedeker, D., Demirtas, H., & Mermelstein, R. J. (2009). A mixed ordinal location scale model for analysis of ecological momentary assessment (EMA) data. *Statistics and Its Interface*, *2*(4), 391.
- Hedges, L. V., & Hedberg, E. C. (2007). Intraclass correlation values for planning group-randomized trials in education. *Educational Evaluation and Policy Analysis*, *29*(1), 60–87. <http://doi.org/10.3102/0162373707299706>
- Ho, A. D. (2008). The problem with “Proficiency”: Limitations of statistics and policy under No Child Left Behind. *Educational Researcher*, *37*(6), 351–360. <http://doi.org/10.3102/0013189X08323842>
- Ho, A. D. (2009). A nonparametric framework for comparing trends and gaps across tests. *Journal of Educational and Behavioral Statistics*, *34*(2), 201–228. <http://doi.org/10.3102/107699860933275>
- Ho, A. D., & Reardon, S. F. (2012). Estimating achievement gaps from test scores reported in ordinal “Proficiency” categories. *Journal of Educational and Behavioral Statistics*, *37*(4), 489–517. <http://doi.org/10.3102/1076998611411918>
- Holland, P. W. (2002). Two measures of change in the gaps between the CDFs of test-score distributions. *Journal of Educational and Behavioral Statistics*, *27*(1), 3–17. <http://doi.org/10.3102/10769986027001003>
- Horowitz, J. L., Sparmann, J. M., & Daganzo, C. F. (1982). An investigation of the accuracy of the clark approximation for the multinomial probit model. *Transportation Science*, *16*(3), 382–401. <http://doi.org/10.1287/trsc.16.3.382>
- Jacob, R. T., Goddard, R. D., & Kim, E. S. (2013). Assessing the use of aggregate data in the evaluation of school-based interventions: Implications for evaluation research and state policy regarding

- public-use data. *Educational Evaluation and Policy Analysis*.
<http://doi.org/10.3102/0162373713485814>
- Jennings, J. (2011). Open Letter to the Member States of PARCC and SBAC. Center on Education Policy. Retrieved from <http://www.cep-dc.org/displayDocument.cfm?DocumentID=359>
- Keane, M. P. (1992). A note on identification in the multinomial probit model. *Journal of Business & Economic Statistics*, *10*(2), 193–200. <http://doi.org/10.1080/07350015.1992.10509898>
- Keele, L., & Park, D. K. (2006, March). *Difficult choices: An evaluation of heterogeneous choice models*. Working Paper. Retrieved from <http://www3.nd.edu/~rwilliam/oglm/ljk-021706.pdf>
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. New York: Routledge.
- McCullagh, P. (1980). Regression models for ordinal data. *Journal of the Royal Statistical Society, Series B (Methodological)*, *42*(2), 109–142.
- Mislevy, R. J., Johnson, E. G., & Muraki, E. (1992). Chapter 3: Scaling procedures in NAEP. *Journal of Educational and Behavioral Statistics*, *17*(2), 131–154.
<http://doi.org/10.3102/10769986017002131>
- Neter, J., Wasserman, W., & Kutner, M. H. (1990). *Applied linear statistical models: Regression, analysis of variance, and experimental designs* (3rd ed.). Homewood, IL: Richard D. Irwin, Inc.
- Reardon, S. F., & Ho, A. D. (2015). Practical issues in estimating achievement gaps from coarsened data. *Journal of Educational and Behavioral Statistics*, *40*(2), 158–189.
<http://doi.org/10.3102/1076998615570944>
- Shear, B. R., Castellano, K. E., & Lockwood, J. R. (2016, April). *Using the Fay-Herriot model to improve inferences from coarsened proficiency data*. Presented at the National Council on Measurement in Education 2016 Annual Meeting, Washington, D.C.
- StataCorp. (2013). *Stata statistical software: Release 13*. College Station, TX: StataCorp LP.
- Tosteson, A. N. A., & Begg, C. B. (1988). A general regression methodology for ROC curve estimation. *Medical Decision Making*, *8*(3), 204–215. <http://doi.org/10.1177/0272989X8800800309>
- U.S. Department of Education. (2015). *State assessments in reading/language arts and mathematics: School year 2012-13 EDFacts Data Documentation*. Washington, D.C. Retrieved from <http://www.ed.gov/edfacts>
- Williams, R. (2009). Using heterogeneous choice models to compare logit and probit coefficients across groups. *Sociological Methods & Research*, *37*(4), 531–559.
<http://doi.org/10.1177/0049124109335735>
- Williams, R. (2010). Fitting heterogeneous choice models with oglm. *The Stata Journal*, *10*(4), 540–567.

Figures and Tables

Table 1. Ratio of Median Estimated Standard Error to Empirical Standard Error and 95% Confidence Interval Coverage for HETOP Estimates by Parameter, Standard Error Formula, Sample Size and Cutscores.

Sample Size	Cutscores	Group Mean		Group Standard Deviation		ICC	
		Ratio	CI	Ratio	CI	Ratio	CI
25	Skewed (05/30/55)	0.923	0.943	0.865	0.892	0.914	0.975
	Mid (20/50/80)	0.944	0.945	0.887	0.906	1.133	0.942
	Wide (05/50/95)	1.608	0.983	1.881	0.995	3.251	0.994
	Many (05/25/50/75/95)	0.965	0.936	0.933	0.927	1.073	0.897
50	Skewed (05/30/55)	0.954	0.946	0.930	0.920	1.084	0.933
	Mid (20/50/80)	0.973	0.948	0.944	0.929	1.084	0.940
	Wide (05/50/95)	1.043	0.957	0.923	0.965	1.476	0.943
	Many (05/25/50/75/95)	0.982	0.943	0.967	0.939	1.044	0.925
100	Skewed (05/30/55)	0.974	0.947	0.964	0.935	1.053	0.924
	Mid (20/50/80)	0.984	0.948	0.972	0.940	1.019	0.933
	Wide (05/50/95)	0.996	0.948	0.957	0.951	1.093	0.929
	Many (05/25/50/75/95)	0.990	0.946	0.983	0.945	1.019	0.932
400	Skewed (05/30/55)	0.995	0.949	0.992	0.946	0.995	0.935
	Mid (20/50/80)	0.999	0.950	0.993	0.948	0.995	0.944
	Wide (05/50/95)	0.998	0.949	0.993	0.949	1.006	0.946
	Many (05/25/50/75/95)	1.000	0.949	0.996	0.948	0.997	0.944

NOTE: ICC = intraclass correlation coefficient. Ratio = ratio of median estimated standard error to empirical standard error. CI = confidence interval coverage rate of an estimated 95% confidence interval.

Table 2. Average, Minimum, and Maximum Correlations between HETOP Estimates and Uncoarsened Score Estimates.

		Means				Standard Deviations			
		(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
Estimate 1:		H4	H20	Orig.	H4	H4	H20	Orig.	H4
Estimate 2:		Orig.	Trans.	Trans.	Trans.	Orig.	Trans.	Trans.	Trans.
NAEP	Average	0.995	1.000	0.999	0.996	0.851	0.995	0.955	0.910
	Minimum	0.988	1.000	0.998	0.992	0.738	0.992	0.929	0.831
	Maximum	0.998	1.000	1.000	0.998	0.923	0.998	0.978	0.967
State	Average	0.973	1.000	0.999	0.973	0.779	0.987	0.932	0.759
	Minimum	0.941	0.999	0.998	0.941	0.667	0.979	0.835	0.626
	Maximum	0.991	1.000	1.000	0.992	0.866	0.995	0.990	0.864

NOTE: H4 = heteroskedastic ordered probit model with 4 proficiency categories as defined by testing program; H20 = heteroskedastic ordered probit model with 20 categories defined by 19 equally spaced cutscores; Orig. = original score scale metric; Trans. = transformed score scale metric.

Figure 1. Average Bias and Aggregate RMSE in Group Standard Deviation Estimates.

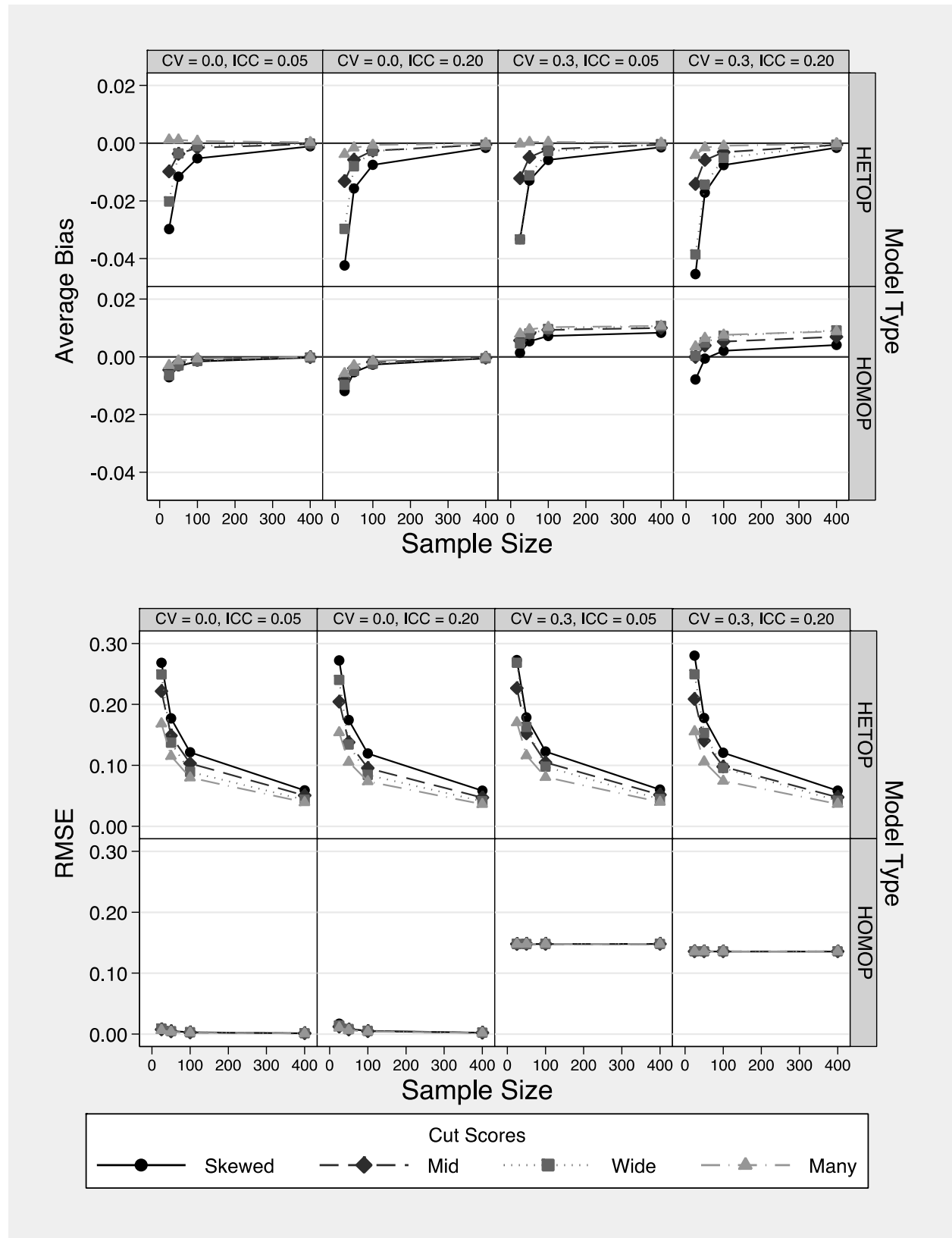


Figure 2. Bias in Standard Deviation Estimates by True Group Mean and Standard Deviation (for ICC=0.20 and CV=0.3 condition)

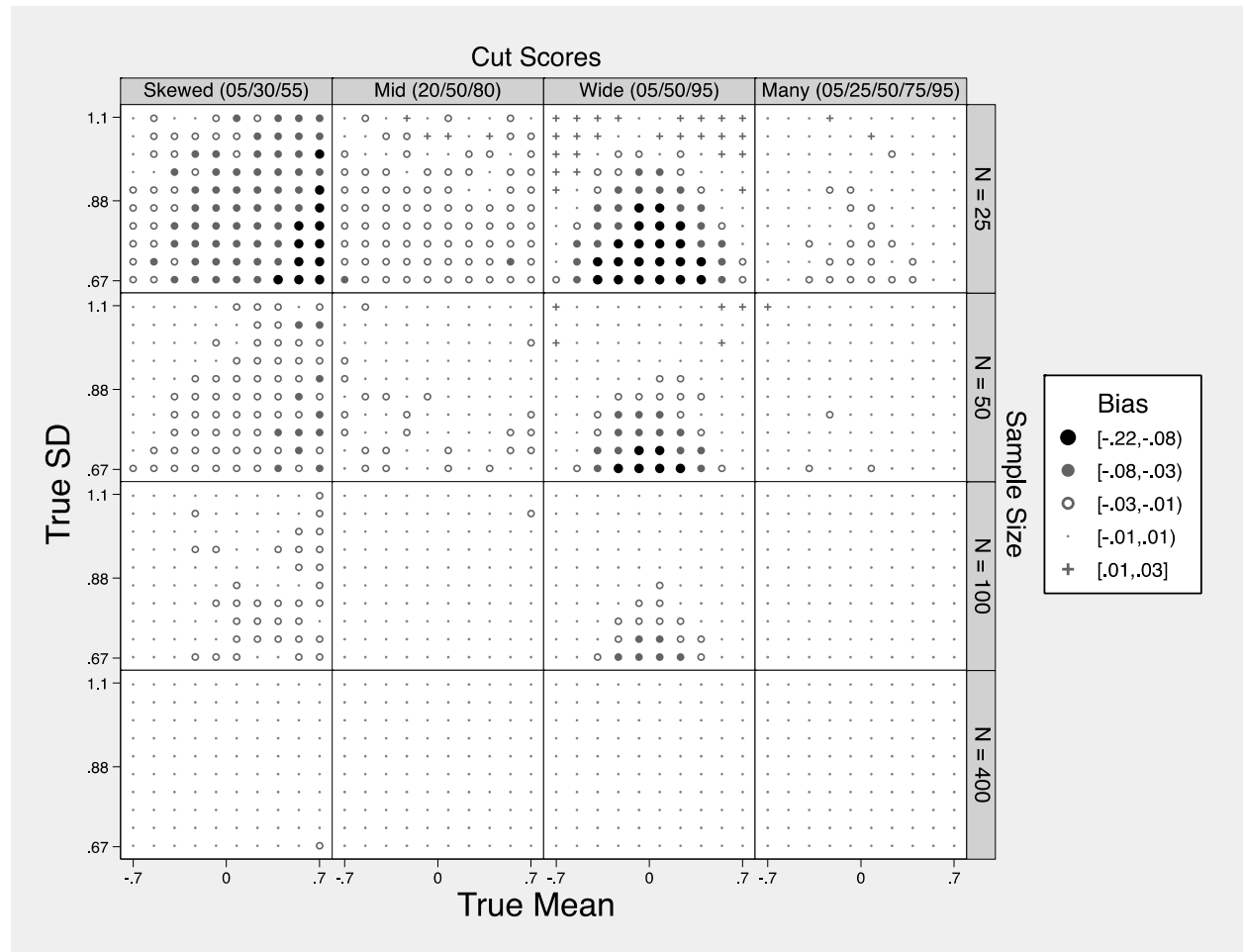


Figure 3. Bias in ICC Estimates

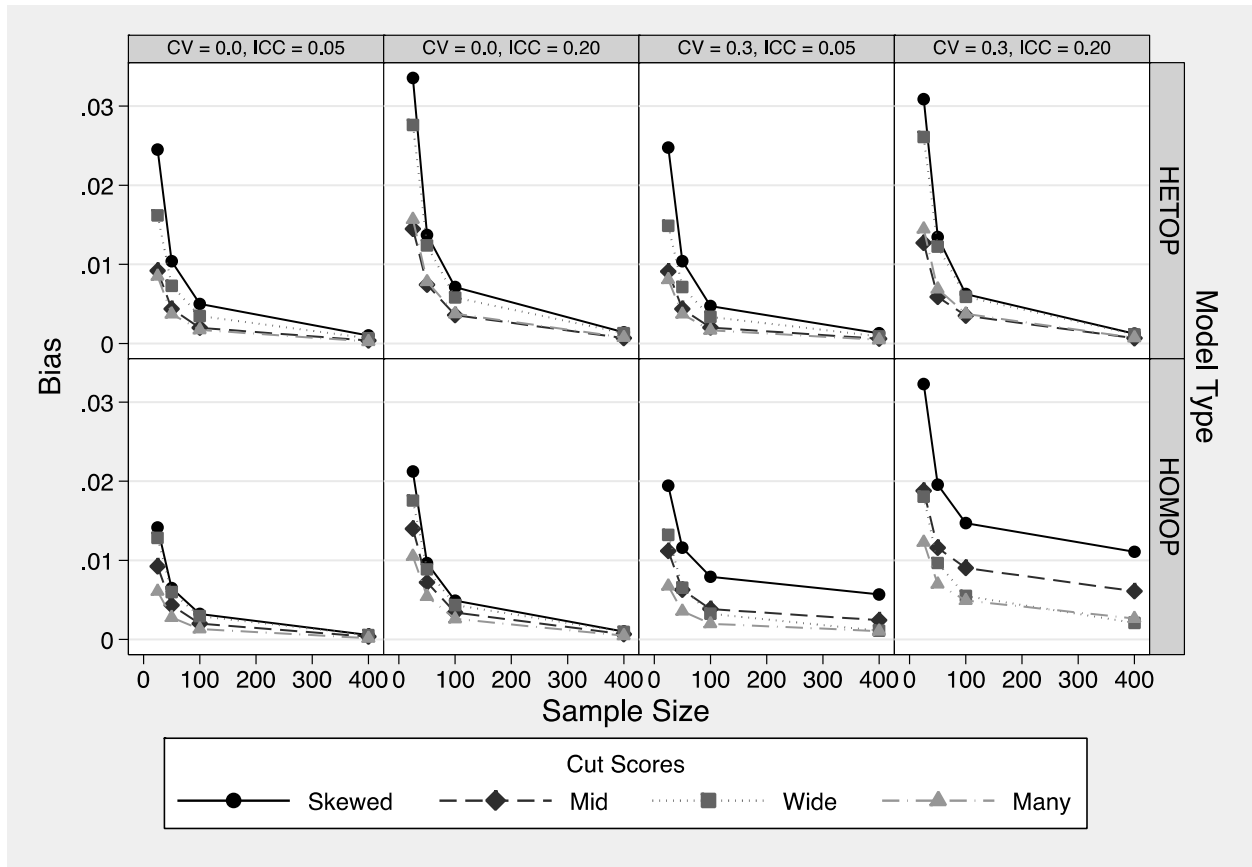


Figure 4. Average Efficiency Ratios of Estimated Means and Standard Deviations Using HETOP Model.

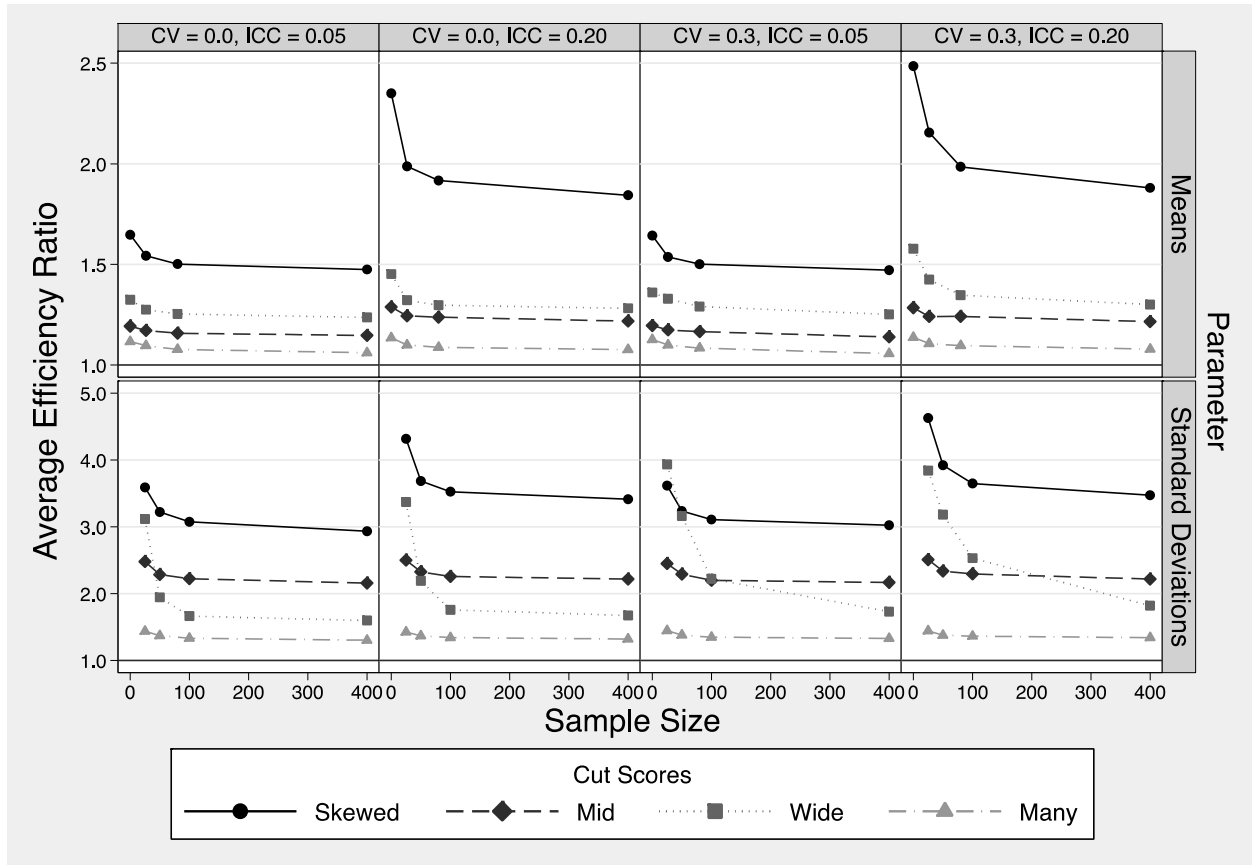
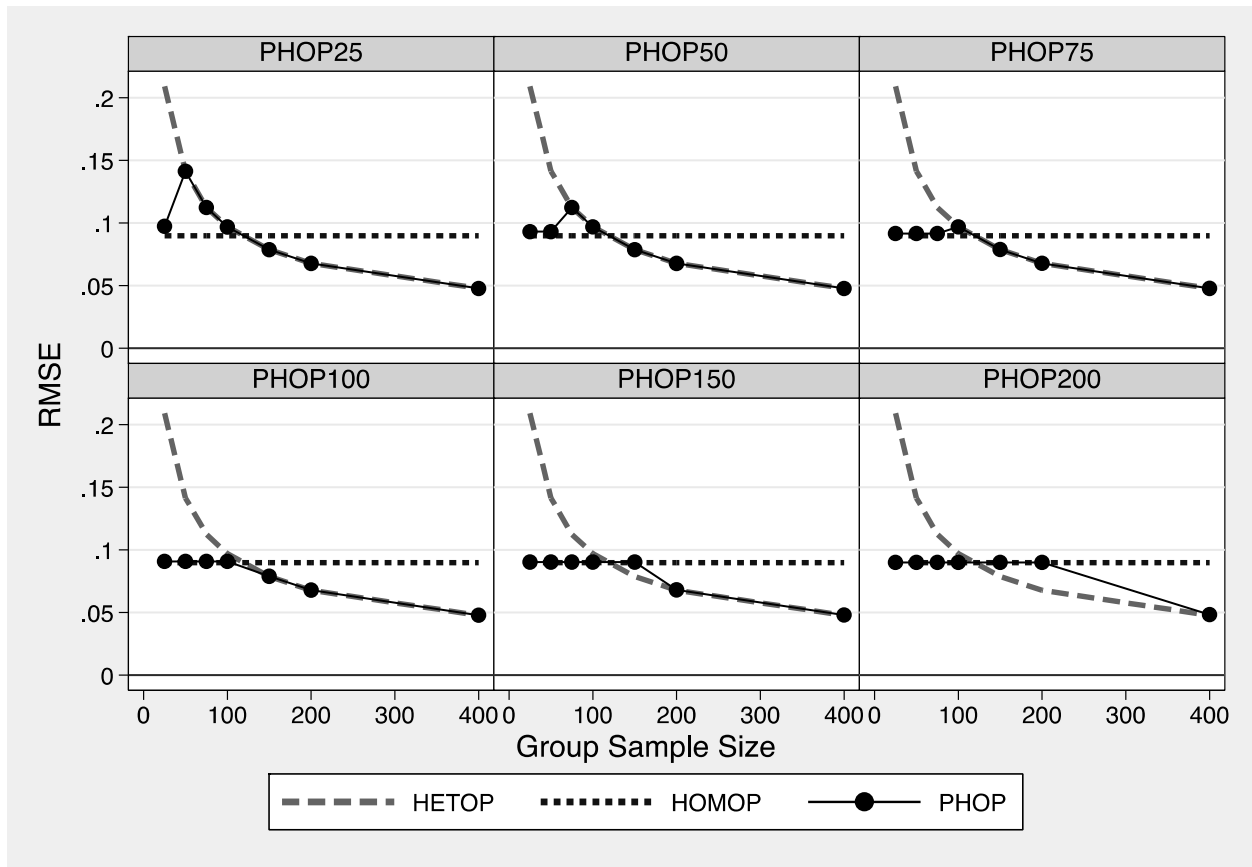
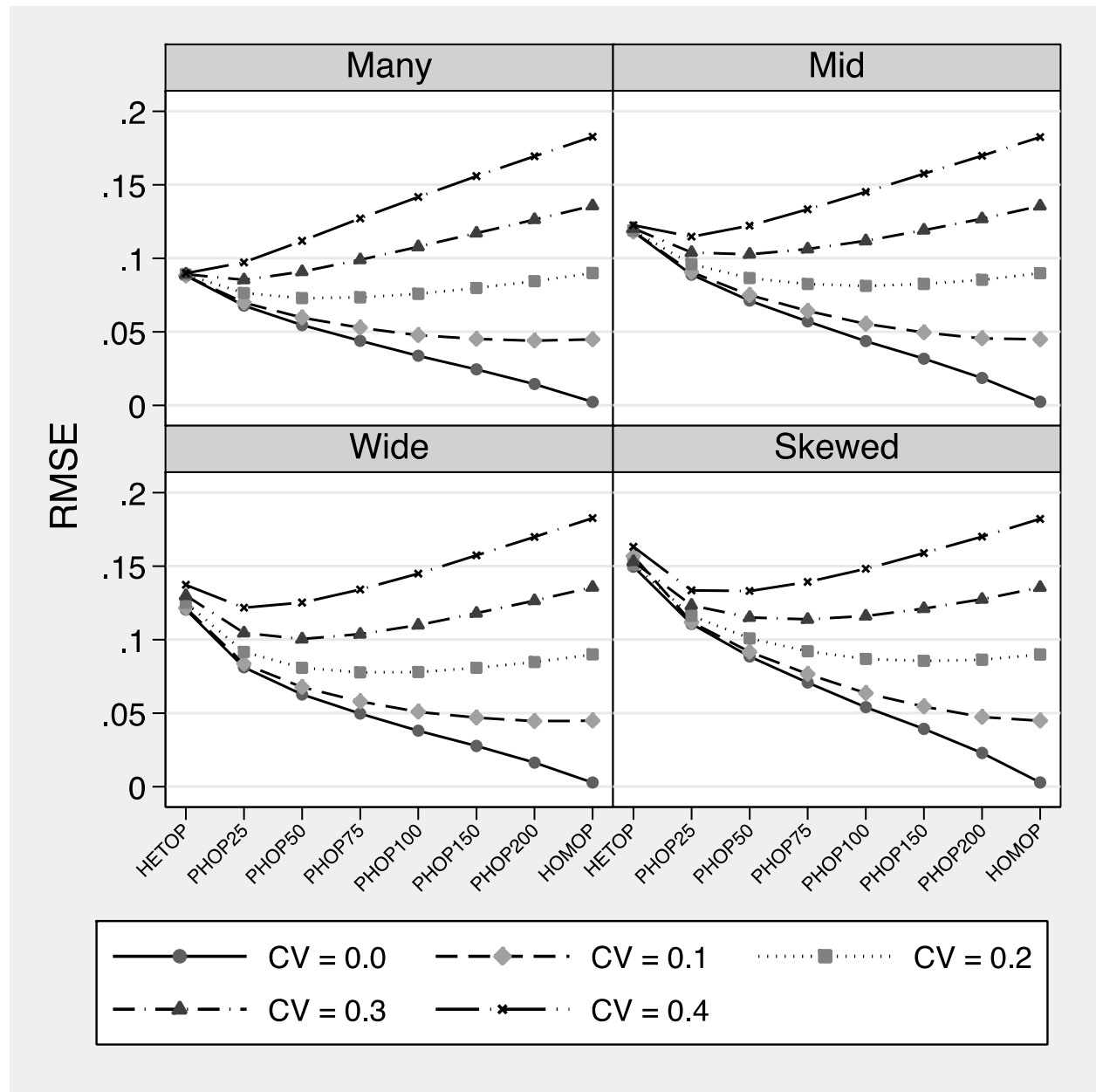


Figure 5. RMSE of Group Standard Deviation Estimates, by Group Sample Size and PHOP Model Type with CV = 0.2



NOTE: HETOP = heteroskedastic ordered probit model; HOMOP = homoskedastic ordered probit model; PHOP = partially heteroskedastic ordered probit model; RMSE = root mean squared error; CV = coefficient of variation. The HETOP and HOMOP lines are included for reference and are constant across all six panels.

Figure 6. Aggregate RMSE of Group Standard Deviation Estimates by Model Type, Cutscore Locations and CV.



Note: CV = coefficient of variation; RMSE = root mean squared error.

Appendix A: Estimating the Total Between- and Within-Group Variances

Given $\hat{\mathbf{M}}'$ and $\hat{\mathbf{\Gamma}}'$, we wish to estimate the within- and between-group variance of y . As noted in the text, we assume throughout this paper that the population consists of a finite number of groups ($g = 1, \dots, G$), all of which are observed. As above, \mathbf{P} is the $1 \times G$ vector of group population proportions (the p_g 's). We observe a sample of size n_g from each group, where n_g may or may not be proportional to p_g . Without loss of generality, we assume the model is fit subject to the constraints that $\mathbf{P}\mathbf{M}' = \mathbf{0}$ and $\mathbf{P}\mathbf{\Gamma}' = \mathbf{0}$. If it is not, we transform the estimate to obtain $\hat{\mathbf{M}}'$ and $\hat{\mathbf{\Sigma}}'$ in this metric, as described in Online Appendix A.

The between- and within group variances are defined as

$$\begin{aligned}\sigma_B'^2 &= \mathbf{P}\mathbf{M}'^2 \\ \sigma_W'^2 &= \mathbf{P}\mathbf{\Sigma}'^2.\end{aligned}\tag{A1}$$

We can compute (biased) estimates of these using their sample analogs, $\mathbf{P}\hat{\mathbf{M}}'^2$ and $\mathbf{P}\hat{\mathbf{\Sigma}}'^2$. Below we derive the expected values of these estimators to assess their bias. We use the results of these derivations to obtain approximately unbiased estimators.

Estimating $\sigma_W'^2$

Let w_g be the error in $\hat{\gamma}_g$: $\hat{\gamma}_g = \gamma_g + \hat{w}_g$. Let $\mathbf{\Omega}'$ be the sampling variance-covariance matrix of the γ_g 's. The diagonal elements of this are the squared sampling variances (the ω_g^2 's). Then

$$\begin{aligned}
E(\mathbf{P}\widehat{\Sigma}^{\circ 2}) &= E\left(\sum_g p_g \widehat{\sigma}_g^2\right) \\
&= \sum_g E(p_g e^{2\widehat{\gamma}_g}) \\
&= \sum_g E\left(p_g e^{2(\gamma_g + \widehat{w}_g)}\right) \\
&= \sum_g E(p_g e^{2\gamma_g} e^{2\widehat{w}_g}) \\
&= \sum_g (p_g \sigma_g^2) E(e^{2\widehat{w}_g}) \\
&\approx \sum_g (p_g \sigma_g^2) E(1 + 2\widehat{w}_g + 2\widehat{w}_g^2) \\
&= \sum_g (p_g \sigma_g^2) (1 + 2\overline{\omega}_g^2) \\
&= \mathbf{P}\Sigma^{\circ 2} \cdot (1 + 2\overline{\omega}_g^2) + 2GCov(p_g \sigma_g^2, \omega_g^2),
\end{aligned} \tag{A2}$$

where $\overline{\omega}_g^2 = \frac{1}{G} \mathbf{1} \cdot \text{vecdiag}(\mathbf{\Omega}')$ is the average sampling variance of the $\widehat{\gamma}_g$'s. Under the assumption that $Cov(p_g \sigma_g^2, \omega_g^2) \approx 0$, we have

$$E(\mathbf{P}\widehat{\Sigma}^{\circ 2}) \approx \mathbf{P}\Sigma^{\circ 2} \cdot (1 + 2\overline{\omega}_g^2). \tag{A3}$$

Therefore, we can compute an approximately unbiased estimate of $\widehat{\sigma}_W'^2$ as

$$\widehat{\sigma}_W'^2 = \frac{\mathbf{P}\widehat{\Sigma}^{\circ 2}}{1 + 2\overline{\omega}_g^2}. \tag{A4}$$

Equation (A4) requires an estimate of $\overline{\omega}_g^2$, the average sampling variance of the $\widehat{\gamma}_g$'s, which can be obtained from the estimated sampling covariance matrix of the $\widehat{\gamma}_g$'s:

$$\widehat{\omega}_g^2 = \frac{1}{G} \mathbf{1} \cdot \text{vecdiag}(\widehat{\mathbf{\Omega}}') \quad (\text{A5})$$

However, $\widehat{\mathbf{\Omega}}'$ is prone to sampling variance (that is, the estimated sampling variances of the γ_g 's themselves have sampling variances). Our simulations show that when n is small, the sampling variance of the elements of $\widehat{\mathbf{\Omega}}'$ can be very large, because the sparse coarsened data provide little information from which to estimate the sampling variances. As a result $E \left[\frac{1}{G} \mathbf{1} \cdot \text{vecdiag}(\widehat{\mathbf{\Omega}}') \right] \gg \frac{1}{G} \mathbf{1} \cdot \text{vecdiag}(\mathbf{\Omega}')$ in such cases.

An alternate method of estimating $\overline{\omega}_g^2$ is to derive an approximate formula based on group sample sizes. To do so, let \hat{u}_g be the error in $\hat{\sigma}_g^2$: $\hat{\sigma}_g^2 = \sigma_g^2 + \hat{u}_g$. If a population variance σ^2 is estimated from a sample of size n (using data that have not been coarsened), the sampling variance of $\hat{\sigma}^2$ is approximately $\frac{2\sigma^4}{n-1}$ (Casella & Berger, 2002; Neter, Wasserman, & Kutner, 1990). Note that, for a normally-distributed variable X with mean 0 and standard deviation s , $\text{var}(X + X^2) \approx s^2 + 2s^4$. We then have

$$\begin{aligned}
\hat{\sigma}_g^2 &= \sigma_g^2 + \hat{u}_g \\
e^{2\hat{\gamma}_g} &= e^{2\gamma_g} + \hat{u}_g \\
e^{2\gamma_g} e^{2\hat{w}_g} &= e^{2\gamma_g} + \hat{u}_g \\
e^{2\gamma_g} (e^{2\hat{w}_g} - 1) &= \hat{u}_g \\
e^{2\gamma_g} (1 + 2\hat{w}_g + 2\hat{w}_g^2 - 1) &\approx \hat{u}_g \\
2e^{2\gamma_g} (\hat{w}_g + \hat{w}_g^2) &= \hat{u}_g \\
\text{var} (2e^{2\gamma_g} (\hat{w}_g + \hat{w}_g^2)) &= \text{var}(\hat{u}_g) \\
4e^{4\gamma_g} \text{var}(\hat{w}_g + \hat{w}_g^2) &= \text{var}(\hat{u}_g) \\
4\sigma_g^4 [\omega_g^2 + 2\omega_g^4] &= \frac{2\sigma_g^4}{n_g - 1} \\
\omega_g^2 + 2\omega_g^4 &= \frac{1}{2(n_g - 1)}
\end{aligned} \tag{A6}$$

Applying the quadratic formula to solve for ω_g^2 yields one positive root:

$$\begin{aligned}
\omega_g^2 &= -\frac{1}{4} + \frac{1}{4} \sqrt{1 + \frac{4}{n_g - 1}} \\
&\approx -\frac{1}{4} + \frac{1}{4} \left(1 + \frac{2}{n_g - 1}\right) \\
&= \frac{1}{2(n_g - 1)},
\end{aligned} \tag{A7}$$

where the approximation holds if n_g is even moderately large.

Given (A7), we have

$$\overline{\omega_g^2} \approx \frac{1}{G} \sum_g \frac{1}{2(n_g - 1)} = \frac{1}{2\bar{n}} \tag{A8}$$

where \tilde{n} is the harmonic mean of $n_g - 1$: $\tilde{n} = \left(\frac{1}{G} \sum_g \frac{1}{n_g - 1} \right)^{-1}$.

Note that (A8) is based on a formula for the sampling variance of a population variance based on uncoarsened data. When the data are coarsened, the sampling variability of $\hat{\sigma}_g'^2$ will certainly be larger than that given by the formula used above $\left(\frac{2\sigma_g'^4}{n_g - 1} \right)$, but the difference may not be large. For example, suppose the true sampling variance of $\sigma_g'^2$ were $\frac{2c_s\sigma_g'^4}{n_g - 1}$, where $c_s \geq 1$; then using the approximation in Equation (A8) in Equation (A4) will inflate our estimate of $\hat{\sigma}_W'^2$ by a factor of $\frac{\tilde{n} + c_s}{\tilde{n} + 1}$. Unless c_s is large in relation to \tilde{n} , the difference will be trivial.

The approximation in Equation (A8) needs to be modified when using either the HOMOP or PHOP (rather than the HETOP) model. When we fit the HOMOP model, the sampling variance in the estimate of the $\hat{\gamma}_g$'s will be smaller, because the estimate is based on the pooled sample of all groups. In this case, $\overline{\omega_g^2}$ might be well-estimated by (A5). Alternately, because the effective sample size for estimating $\overline{\omega_g^2}$ is N , and we lose a degree of freedom in estimating each group's mean, (A8) can be replaced by

$$\overline{\omega_g^2} \approx \frac{1}{2(N - G)}. \quad (\text{A9})$$

In the PHOP model, the average sampling variance of the γ_g 's can be approximated as

$$\overline{\omega_g^2} \approx \frac{1}{G} \sum_g \frac{1}{2(\check{n}_g - 1)}, \quad (\text{A10})$$

where $\check{n}_g = n_g$ if group g 's standard deviation is not constrained, and $\check{n}_g = \sum_{g \in C} (n_g - 1)$ if group $g \in C$, where C is the set of constrained groups.

When estimating $\hat{\sigma}_W'^2$, we substitute either Equation (A8), (A9), or (A10) into Equation (A4) depending upon which model was fit.

Estimating $\sigma_B'^2$

To compute the expected value of $\mathbf{P}\hat{\mathbf{M}}'^{\circ 2}$, first note that estimating the variance of the group means involves error in the overall mean and the individual group means. The estimate of each group's mean has two sources of error in it: $\hat{\mu}'_g = \mu_g - \hat{u} + \hat{e}_g$, where $\hat{u} = \sum p_g \hat{e}_g$ and $\hat{e}_g = \hat{\mu}'_g - \mu_g$. Then

$$\begin{aligned}
 E[\mathbf{P}\hat{\mathbf{M}}'^{\circ 2}] &= E[\mathbf{P}(\mathbf{M} - \hat{\mathbf{u}}' + \hat{\mathbf{e}}')^{\circ 2}] \\
 &= E[\mathbf{P}\mathbf{M}^{\circ 2} - 2\mathbf{P}(\hat{\mathbf{u}}' \circ \hat{\mathbf{e}}') + \mathbf{P}\hat{\mathbf{u}}'^{\circ 2} + \mathbf{P}\hat{\mathbf{e}}'^{\circ 2}] \\
 &= \mathbf{P}\mathbf{M}^{\circ 2} - 2\mathbf{P}E[\hat{\mathbf{u}}' \circ \hat{\mathbf{e}}'] + \mathbf{P}E[\hat{\mathbf{u}}'^{\circ 2}] + \mathbf{P}E[\hat{\mathbf{e}}'^{\circ 2}] \\
 &= \mathbf{P}\mathbf{M}^{\circ 2} - \mathbf{P}\mathbf{V}'\mathbf{P}^t + \mathbf{P} \cdot \mathit{vecdiag}(\mathbf{V}')
 \end{aligned} \tag{A11}$$

where $\mathit{vecdiag}(\mathbf{V}')$ is the $G \times 1$ matrix of sampling variances of the means (the diagonal of \mathbf{V}'). So we can compute an unbiased estimate of $\hat{\sigma}_B'^2$ as

$$\hat{\sigma}_B'^2 = \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \mathbf{P}\mathbf{V}'\mathbf{P}^t - \mathbf{P} \cdot \mathit{vecdiag}(\mathbf{V}'). \tag{A12}$$

Equation (A12) requires an estimate of \mathbf{V}' , the variance-covariance matrix of the vector of estimated group means, $\hat{\mathbf{M}}'$. One estimate of this is the estimated matrix $\hat{\mathbf{V}}'$. However, like $\hat{\mathbf{\Omega}}'$ above, $\hat{\mathbf{V}}'$ is prone to sampling variance (that is, the estimated sampling variances of the $\hat{\mu}'_g$'s themselves have sampling variances). Our simulations show that when n is small, the sampling variance of the elements of $\hat{\mathbf{V}}'$ can be very large, because the sparse coarsened data provide little information from which to estimate the sampling variances. As a result $E[\mathbf{P} \cdot \mathit{vecdiag}(\hat{\mathbf{V}}')] \gg \mathbf{P} \cdot \mathit{vecdiag}(\mathbf{V}')$ in such cases.

An alternate method of estimating $\hat{\mathbf{V}}'$ is to derive an approximate formula for its diagonal elements based on group sample sizes. We begin by assuming that the off-diagonal elements of \mathbf{V}' are approximately 0 (they will not be exactly zero, because the estimated means are dependent on one another, since all are estimated simultaneously and constrained to satisfy $\mathbf{P}\hat{\mathbf{M}}' = \mathbf{0}$, but they will be close to zero when G and n are moderately large). We then assume the sampling variance of μ is given by the standard formula (based on uncoarsened data) for the sampling variance of a mean: $\mathit{var}(\hat{\mu}) = \frac{\sigma^2}{n}$. Then

the diagonal elements of \mathbf{V}' will be $v_{gg} = \frac{\sigma_g'^2}{n_g}$. Substituting this matrix into (A12), we get

$$\begin{aligned}
\hat{\sigma}_B'^2 &= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \mathbf{P}\mathbf{V}'\mathbf{P}^t - \mathbf{P} \cdot \text{vecdiag}(\mathbf{V}') \\
&= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + (\mathbf{P}^{\circ 2} - \mathbf{P}) \cdot \text{vecdiag}(\mathbf{V}') \\
&= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \sum_g (p_g^2 - p_g) \frac{\sigma_g'^2}{n_g} \\
&= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + (\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P}))\boldsymbol{\Sigma}'^{\circ 2} \\
&= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \frac{(\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P}))\hat{\boldsymbol{\Sigma}}'^{\circ 2}}{1 + 2\widehat{\omega}_g^2} \\
&= \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \frac{(\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P}))\hat{\boldsymbol{\Sigma}}'^{\circ 2}}{1 + 2\widehat{\omega}_g^2},
\end{aligned} \tag{A13}$$

where we substitute in the approximations of $\widehat{\omega}_g^2$ from Equations (A8), (A9), or (A10) as appropriate.

Again, if the sampling variance of the $\hat{\mu}'_g$ estimates is greater than they would be if the data were not coarsened, then it may be more appropriate to substitute $v_{gg} = c_m \frac{\sigma_g'^2}{n_g}$ into (A12) above, where $c_m \geq 1$ is a constant. Then if we also use c_s as above in the formula estimating $\boldsymbol{\Sigma}'^{\circ 2}$, (A13) becomes:

$$\hat{\sigma}_B'^2 = \mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \frac{c_m}{1 + c_s 2\widehat{\omega}_g^2} (\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P}))\hat{\boldsymbol{\Sigma}}'^{\circ 2}. \tag{A14}$$

Given that $(\mathbf{n}^{\circ -1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P}))\hat{\boldsymbol{\Sigma}}'^{\circ 2}$ will be small when the elements of \mathbf{n} are modestly large, however, setting $c_m = 1$ has very little effect of the estimate of $\hat{\sigma}_B'^2$.

Estimating the population standard deviation, σ'

Given estimates of $\sigma_W'^2$ and $\sigma_B'^2$ from (A4) and (A13), we compute

$$\hat{\sigma}' = (\sigma_W'^2 + \sigma_B'^2)^{\frac{1}{2}} \tag{A15}$$

$$\begin{aligned}
&= \left(\mathbf{P}\widehat{\mathbf{M}}^{\prime\circ 2} + \frac{\mathbf{n}^{\circ-1} \circ (\mathbf{P}^{\circ 2} - \mathbf{P})\widehat{\boldsymbol{\Sigma}}^{\prime\circ 2}}{1 + 2\widehat{\omega}_g^2} + \frac{\mathbf{P}\widehat{\boldsymbol{\Sigma}}^{\prime\circ 2}}{1 + 2\widehat{\omega}_g^2} \right)^{\frac{1}{2}} \\
&= \left(\mathbf{P}\widehat{\mathbf{M}}^{\prime\circ 2} + \frac{(\mathbf{n}^{\circ-1} \circ (\mathbf{P} + \mathbf{n} - \mathbf{1}) \circ \mathbf{P})\widehat{\boldsymbol{\Sigma}}^{\prime\circ 2}}{1 + 2\widehat{\omega}_g^2} \right)^{\frac{1}{2}} \\
&= (\mathbf{P}\widehat{\mathbf{M}}^{\prime\circ 2} + \mathbf{Q}\widehat{\boldsymbol{\Sigma}}^{\prime\circ 2})^{\frac{1}{2}},
\end{aligned}$$

where

$$\mathbf{Q} = \frac{(\mathbf{n}^{\circ-1} \circ (\mathbf{P} + \mathbf{n} - \mathbf{1}) \circ \mathbf{P})}{1 + 2\widehat{\omega}_g^2}, \tag{A16}$$

And we again substitute one of the approximations from Equations (A8), (A9), or (A10) for $\overline{\omega}_g^2$ depending upon whether a HETOP, HOMOP, or PHOP model is fit.

Appendix B: Computation of standard errors of $\widehat{\mathbf{M}}^*$ and $\widehat{\boldsymbol{\Sigma}}^*$

Once we have constructed $\widehat{\mathbf{M}}^*$ and $\widehat{\boldsymbol{\Sigma}}^*$ via Equation (6), we must estimate the covariance matrices $\mathbf{V}^* = \text{Cov}(\widehat{\mathbf{M}}^*, \widehat{\mathbf{M}}^*)$, $\mathbf{Z}^* = \text{Cov}(\widehat{\mathbf{M}}^*, \widehat{\boldsymbol{\Sigma}}^*)$, and $\mathbf{W}^* = \text{Cov}(\widehat{\boldsymbol{\Sigma}}^*, \widehat{\boldsymbol{\Sigma}}^*)$, from which we can obtain standard errors for the parameters of interest in the model.

Assuming that $E[\widehat{\mathbf{M}}'] = \mathbf{M}'$ and $E[\widehat{\boldsymbol{\Sigma}}'] = \boldsymbol{\sigma}'$,¹⁰ the g, h element \mathbf{V}^* is

$$\begin{aligned}
 v_{gh}^* &= \text{Cov}(\hat{\mu}_g^*, \hat{\mu}_h^*) \\
 &= \text{Cov}\left(\frac{\hat{\mu}'_g}{\hat{\sigma}'}, \frac{\hat{\mu}'_h}{\hat{\sigma}'}\right) \\
 &\approx \frac{1}{\sigma'^2} v'_{gh} - \frac{\mu'_g}{\sigma'^3} \text{Cov}(\hat{\sigma}', \hat{\mu}'_h) - \frac{\mu'_h}{\sigma'^3} \text{Cov}(\hat{\mu}'_g, \hat{\sigma}') + \hat{\mu}'_g \hat{\mu}'_h \text{Var}\left(\frac{1}{\hat{\sigma}'}\right) \\
 &\approx \frac{1}{\sigma'^2} [v'_{gh} - \mu_g^* \text{Cov}(\hat{\sigma}', \hat{\mu}'_h) - \mu_h^* \text{Cov}(\hat{\mu}'_g, \hat{\sigma}') + \mu_g^* \mu_h^* \text{Var}(\hat{\sigma}')].
 \end{aligned} \tag{B1}$$

Now let \mathbf{I}_h denote the h^{th} column of the $G \times G$ identity matrix. Then define¹¹

$$\begin{aligned}
 r_h &= \text{Cov}(\hat{\sigma}', \hat{\mu}'_h) \\
 &= \frac{1}{2\sigma'} \text{Cov}(\hat{\sigma}'^2, \hat{\mu}'_h) \\
 &= \frac{1}{2\sigma'} \text{Cov}(\mathbf{P}\widehat{\mathbf{M}}'^{\circ 2} + \mathbf{Q}\widehat{\boldsymbol{\Sigma}}'^{\circ 2}, \hat{\mu}'_h) \\
 &= \frac{1}{\sigma'} \mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{V}' \mathbf{I}_h + \frac{1}{\sigma'} \mathbf{Q}[\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Z}'^t \mathbf{I}_h.
 \end{aligned} \tag{B2}$$

Then define the $1 \times G$ vector \mathbf{R} , with elements r_h , as

$$\begin{aligned}
 \mathbf{R} &= \frac{1}{\sigma'} \mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{V}' + \frac{1}{\sigma'} \mathbf{Q}[\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Z}'^t \\
 &= \mathbf{P}[\text{diag}(\mathbf{M}^*)] \mathbf{V}' + \mathbf{Q}[\text{diag}(\boldsymbol{\Sigma}^*)] \mathbf{Z}'^t.
 \end{aligned} \tag{B3}$$

¹⁰ Even under the assumption that the HETOP estimator provides unbiased estimates of \mathbf{M}' and $\boldsymbol{\Sigma}'$, the assumption that $E[\widehat{\boldsymbol{\Sigma}}'] = \boldsymbol{\sigma}'$ is not strictly valid, given nonlinearities in Equations (8) and (9), but is a good approximation in practice.

¹¹ Note that \mathbf{Q} in these formulas depends upon whether a HETOP, HOMOP, or PHOP model is being used, as defined in (A16).

Now we have

$$\mathbf{V}^* \approx \frac{1}{\sigma'^2} [\mathbf{V}' - (\mathbf{M}^* \mathbf{R} + \mathbf{R}^t \mathbf{M}^{*t}) + \mathbf{M}^* \mathbf{M}^{*t} \text{Var}(\hat{\sigma}')]. \quad (\text{B4})$$

Similarly, assuming that $E[\hat{\Sigma}'] = \Sigma'$ and $E[\hat{\sigma}'] = \sigma'$, the g, h element of the covariance matrix

\mathbf{W}^* of the $\hat{\sigma}_g^*$'s is

$$\begin{aligned} w_{gh}^* &= \text{Cov}(\hat{\sigma}_g^*, \hat{\sigma}_h^*) \\ &= \text{Cov}\left(\frac{\hat{\sigma}_g'}{\hat{\sigma}'}, \frac{\hat{\sigma}_h'}{\hat{\sigma}'}\right) \\ &\approx \frac{1}{\sigma'^2} w'_{gh} - \frac{\sigma_g'}{\sigma'^3} \text{Cov}(\hat{\sigma}', \hat{\sigma}_h') - \frac{\sigma_h'}{\sigma'^3} \text{Cov}(\hat{\sigma}_g', \hat{\sigma}') + \hat{\sigma}_g' \hat{\sigma}_h' \text{Var}\left(\frac{1}{\sigma'}\right) \\ &\approx \frac{1}{\sigma'^2} [w'_{gh} - \sigma_g^* \text{Cov}(\hat{\sigma}', \hat{\sigma}_h') - \sigma_h^* \text{Cov}(\hat{\sigma}_g', \hat{\sigma}') + \sigma_g^* \sigma_h^* \text{Var}(\hat{\sigma}')]. \end{aligned} \quad (\text{B5})$$

Define

$$\begin{aligned} t_h &= \text{Cov}(\hat{\sigma}', \hat{\sigma}_h') \\ &= \frac{1}{2\sigma'} \text{Cov}(\hat{\sigma}'^2, \hat{\sigma}_h') \\ &= \frac{1}{2\sigma'} \text{Cov}(\mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \mathbf{Q}\hat{\Sigma}'^{\circ 2}, \hat{\sigma}_h') \\ &= \frac{1}{\sigma'} \mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{Z}' \mathbf{I}_h + \frac{1}{\sigma'} \mathbf{Q}[\text{diag}(\Sigma')] \mathbf{W}' \mathbf{I}_h. \end{aligned} \quad (\text{B6})$$

Then define the $1 \times G$ vector \mathbf{T} , with elements t_h , as

$$\begin{aligned} \mathbf{T} &= \frac{1}{\sigma'} \mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{Z}' + \frac{1}{\sigma'} \mathbf{Q}[\text{diag}(\Sigma')] \mathbf{W}' \\ &= \mathbf{P}[\text{diag}(\mathbf{M}^*)] \mathbf{Z}' + \mathbf{Q}[\text{diag}(\Sigma^*)] \mathbf{W}'. \end{aligned} \quad (\text{B7})$$

We then have

$$\mathbf{W}^* \approx \frac{1}{\sigma'^2} [\mathbf{W}' - (\Sigma^* \mathbf{T} + \mathbf{T}^t \Sigma^{*t}) + \Sigma^* \Sigma^{*t} \text{Var}(\hat{\sigma}')]. \quad (\text{B8})$$

Finally, the element z_{gh}^* of the matrix \mathbf{Z}^* is

$$\begin{aligned}
z_{gh}^* &= \text{Cov}(\hat{\mu}_g^*, \hat{\sigma}_h^*) \\
&= \text{Cov}\left(\frac{\hat{\mu}'_g}{\hat{\sigma}'^t}, \frac{\hat{\sigma}'_h}{\hat{\sigma}'^t}\right) \\
&\approx \frac{1}{\sigma'^2} z'_{gh} - \frac{\mu'_g}{\sigma'^3} \text{Cov}(\hat{\sigma}', \hat{\sigma}'_h) - \frac{\sigma'_h}{\sigma'^3} \text{Cov}(\hat{\mu}'_g, \hat{\sigma}') + \hat{\mu}'_g \hat{\sigma}'_h \text{Var}\left(\frac{1}{\sigma'}\right) \\
&\approx \frac{1}{\sigma'^2} [z'_{gh} - \mu_g^* \text{Cov}(\hat{\sigma}', \hat{\sigma}'_h) - \sigma_h^* \text{Cov}(\hat{\mu}'_g, \hat{\sigma}') + \mu_g^* \sigma_h^* \text{Var}(\hat{\sigma}')] \\
&= \frac{1}{\sigma'^2} [z'_{gh} - \mu_g^* t_h - \sigma_h^* r_g + \mu_g^* \sigma_h^* \text{Var}(\hat{\sigma}')].
\end{aligned} \tag{B9}$$

We then have

$$\mathbf{Z}^* \approx \frac{1}{\sigma'^2} [\mathbf{Z}' - (\mathbf{M}^* \mathbf{T} + \mathbf{R}^t \boldsymbol{\Sigma}^{*t}) + \mathbf{M}^* \boldsymbol{\Sigma}^{*t} \text{Var}(\hat{\sigma}')]. \tag{B10}$$

Expressions (B4), (B8) and (B10) require $\text{Var}(\hat{\sigma}')$. Note first, that we can derive¹² the sampling variance of $\hat{\sigma}'^2$ as

$$\begin{aligned}
\text{Var}(\hat{\sigma}'^2) &= \text{Var}(\mathbf{P}\hat{\mathbf{M}}'^{\circ 2} + \mathbf{Q}\hat{\boldsymbol{\Sigma}}'^{\circ 2}) \\
&= 4\mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{V}' [\text{diag}(\mathbf{M}')] \mathbf{P}^t + 4\mathbf{Q}[\text{diag}(\boldsymbol{\Sigma}')] \mathbf{W}' [\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Q}^t \\
&\quad + 8\mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{Z}' [\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Q}^t.
\end{aligned} \tag{B11}$$

Then, by the Delta method,

$$\begin{aligned}
\text{Var}(\hat{\sigma}') &\approx \frac{1}{4\sigma'^2} \text{Var}(\hat{\sigma}'^2) \\
&\approx \frac{1}{\sigma'^2} [\mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{V}' [\text{diag}(\mathbf{M}')] \mathbf{P}^t + \mathbf{Q}[\text{diag}(\boldsymbol{\Sigma}')] \mathbf{W}' [\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Q}^t \\
&\quad + 2\mathbf{P}[\text{diag}(\mathbf{M}')] \mathbf{Z}' [\text{diag}(\boldsymbol{\Sigma}')] \mathbf{Q}^t].
\end{aligned} \tag{B12}$$

¹² Note that if \mathbf{A} and \mathbf{B} are $1 \times G$ scalar vectors, \mathbf{D} and \mathbf{E} are $G \times G$ scalar matrices; and $\hat{\mathbf{X}}$ and $\hat{\mathbf{Y}}$ are $G \times 1$ column vectors of random variables, then

$$\text{Cov}\left(\mathbf{A}[\mathbf{D}\hat{\mathbf{X}}]^{\circ 2}, \mathbf{B}[\mathbf{E}\hat{\mathbf{Y}}]^{\circ 2}\right) \approx 4\mathbf{A}[\text{diag}(\mathbf{X})][\mathbf{D}^t \mathbf{D}]\mathbf{C}[\mathbf{E}\mathbf{E}^t][\text{diag}(\mathbf{Y})]\mathbf{B}^t,$$

where \mathbf{C} is the $G \times G$ covariance matrix of $\hat{\mathbf{X}}$ and $\hat{\mathbf{Y}}$.

We substitute (B12) into (B4), (B8), and (B10) to obtain expressions for \mathbf{V}^* , \mathbf{W}^* , and \mathbf{Z}^* . To estimate \mathbf{V}^* , \mathbf{W}^* , and \mathbf{Z}^* , we replace the relevant terms in the resulting expressions by their estimated values.

Online Appendix A: Matrix Conversions

The maximization of Equation (4) yields estimates of \mathbf{M}' , $\mathbf{\Gamma}'$, and \mathbf{C}' , subject to a set of linear constraints that define their metric. Although we suggest the constraints $\mathbf{P}\mathbf{M}' = \mathbf{0}$ and $\mathbf{P}\mathbf{\Gamma}' = \mathbf{0}$, there may be cases where the model does not converge under these constraints, or where other constraints are preferable. In such cases, before using Equation (6) to estimate \mathbf{M}^* and $\mathbf{\Sigma}^*$, it is necessary to convert the parameter estimates and their covariance matrix to the metric in which $\mathbf{P}\widehat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\widehat{\mathbf{\Gamma}}' = \mathbf{0}$. We do this as follows.

First, we will use a “double prime” subscript to denote estimates subject to some arbitrary set of constraints. So $\widehat{\mathbf{M}}''$, $\widehat{\mathbf{\Gamma}}''$, and $\widehat{\mathbf{C}}''$ are the estimated parameters of the model. The covariance matrices of $\widehat{\mathbf{M}}''$ and $\widehat{\mathbf{\Gamma}}''$ are denoted $\mathbf{V}'' = Cov(\widehat{\mathbf{M}}'', \widehat{\mathbf{M}}'')$, $\mathbf{\Lambda}'' = Cov(\widehat{\mathbf{M}}'', \widehat{\mathbf{\Gamma}}'')$, and $\mathbf{\Omega}'' = Cov(\widehat{\mathbf{\Gamma}}'', \widehat{\mathbf{\Gamma}}'')$. We want to convert these to the metric (denoted by a “single prime” superscript) in which $\mathbf{P}\widehat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\widehat{\mathbf{\Gamma}}' = \mathbf{0}$. To do so, we define $\mathbf{\Pi} = \mathbf{I} - \mathbf{1}^t\mathbf{P}$ and $\kappa = 1 + \frac{1}{2}var(\mathbf{P}\widehat{\mathbf{\Gamma}}'') = 1 + \frac{1}{2}\mathbf{P}\widehat{\mathbf{\Omega}}''\mathbf{P}^t$. We then construct:

$$\begin{aligned}\widehat{\mathbf{M}}' &= e^{(-\mathbf{P}\widehat{\mathbf{\Gamma}}'')}[\mathbf{\Pi}\widehat{\mathbf{M}}''] \\ \widehat{\mathbf{\Gamma}}' &= \mathbf{\Pi}\widehat{\mathbf{\Gamma}}'' \\ \widehat{\mathbf{C}}' &= e^{(-\mathbf{P}\widehat{\mathbf{\Gamma}}'')}[\widehat{\mathbf{C}}'' - \mathbf{1}^t\mathbf{P}\widehat{\mathbf{M}}'']\end{aligned}\tag{A1}$$

Note that the vector of group standard deviations in the “double prime” metric will be $\widehat{\mathbf{\Sigma}}'' = \exp(\widehat{\mathbf{\Gamma}}'')$ (where $\exp(\mathbf{X})$ denotes the matrix whose elements are the exponentiated values of the corresponding elements of the matrix \mathbf{X}). The standard deviations in the new “single prime” metric will be $\widehat{\mathbf{\Sigma}}' = \exp(\widehat{\mathbf{\Gamma}}') = \exp(\mathbf{\Pi}\widehat{\mathbf{\Gamma}}'') = e^{(-\mathbf{P}\widehat{\mathbf{\Gamma}}'')} \widehat{\mathbf{\Sigma}}''$. So the transformations in (A1) ensure that the group means and standard deviations, as well as the estimated thresholds in \mathbf{C} are all scaled by the same factor, $e^{(-\mathbf{P}\widehat{\mathbf{\Gamma}}'')}$. It is straightforward to show that the transformations also ensure that $\mathbf{P}\widehat{\mathbf{M}}' = \mathbf{0}$ and $\mathbf{P}\widehat{\mathbf{\Gamma}}' = \mathbf{0}$ (because $\mathbf{P}\mathbf{\Pi} = \mathbf{P}\mathbf{I} - \mathbf{P}\mathbf{1}^t\mathbf{P} = \mathbf{P} - \mathbf{P} = \mathbf{0}$). Thus, these transformations represent a linear transformation of the

“double prime” metric that meets the constraints of the “single prime metric.”¹³

Once $\hat{\mathbf{M}}''$ and $\hat{\mathbf{\Gamma}}''$ are transformed, we must also determine the corresponding transformations of the \mathbf{V}'' , $\mathbf{\Lambda}''$, and $\mathbf{\Omega}''$ covariance matrices. Let $\mathbf{\Pi}_g$ denote the g^{th} row of $\mathbf{\Pi}$. Note that $\mathbf{\Pi}_g \hat{\mathbf{M}}'' = \hat{\mu}'_g \cdot e^{(\mathbf{P}\hat{\mathbf{\Gamma}}'')}$.

First, we want to compute the g, h element of \mathbf{V}' :

$$v'_{gh} = cov(\hat{\mu}'_g, \hat{\mu}'_h) = cov\left(\frac{\mathbf{\Pi}_g \hat{\mathbf{M}}''}{e^{(\mathbf{P}\hat{\mathbf{\Gamma}}'')}} , \frac{\mathbf{\Pi}_h \hat{\mathbf{M}}''}{e^{(\mathbf{P}\hat{\mathbf{\Gamma}}'')}}\right). \quad (\text{A2})$$

We assume that $E[\hat{\mathbf{M}}''] = \mathbf{M}''$, $E[\hat{\mathbf{\Gamma}}''] = \mathbf{\Gamma}''$, and use the Taylor series approximations below:

$$\begin{aligned} E[e^{(\mathbf{P}\hat{\mathbf{\Gamma}}'')}] &\approx E\left[e^{(E[\mathbf{P}\hat{\mathbf{\Gamma}}''])} + (\mathbf{P}\hat{\mathbf{\Gamma}}'' - E[\mathbf{P}\hat{\mathbf{\Gamma}}''])e^{(E[\mathbf{P}\hat{\mathbf{\Gamma}}''])} + \frac{1}{2}(\mathbf{P}\hat{\mathbf{\Gamma}}'' - E[\mathbf{P}\hat{\mathbf{\Gamma}}''])^2 e^{(E[\mathbf{P}\hat{\mathbf{\Gamma}}''])} + \dots\right] \\ &\approx e^{(E[\mathbf{P}\hat{\mathbf{\Gamma}}''])} \left[1 + \frac{1}{2}var(\mathbf{P}\hat{\mathbf{\Gamma}}'')\right] \\ &\approx \kappa e^{(\mathbf{P}\mathbf{\Gamma}''}); \\ E[e^{(-\mathbf{P}\hat{\mathbf{\Gamma}}'')}] &\approx E\left[e^{(E[-\mathbf{P}\hat{\mathbf{\Gamma}}''])} - (\mathbf{P}\hat{\mathbf{\Gamma}}'' - E[\mathbf{P}\hat{\mathbf{\Gamma}}''])e^{(E[-\mathbf{P}\hat{\mathbf{\Gamma}}''])} + \frac{1}{2}(\mathbf{P}\hat{\mathbf{\Gamma}}'' - E[\mathbf{P}\hat{\mathbf{\Gamma}}''])^2 e^{(E[-\mathbf{P}\hat{\mathbf{\Gamma}}''])} + \dots\right] \\ &\approx e^{(E[-\mathbf{P}\hat{\mathbf{\Gamma}}''])} \left[1 + \frac{1}{2}var(\mathbf{P}\hat{\mathbf{\Gamma}}'')\right] \\ &\approx \kappa e^{(-\mathbf{P}\mathbf{\Gamma}''),} \end{aligned} \quad (\text{A3})$$

where $\kappa = 1 + \frac{1}{2}var(\mathbf{P}\hat{\mathbf{\Gamma}}'') = 1 + \frac{1}{2}\mathbf{P}\mathbf{\Omega}''\mathbf{P}^t$. For compactness of notation, denote $\hat{a} = \mathbf{\Pi}_g \hat{\mathbf{M}}''$; $\hat{b} =$

$\mathbf{\Pi}_h \hat{\mathbf{M}}''$; $\hat{c} = e^{(\mathbf{P}\hat{\mathbf{\Gamma}}'')}$. Then

¹³ Note that if \mathbf{M}'' , $\mathbf{\Gamma}''$, and \mathbf{C}'' are estimated using the constraints that $\mathbf{P}\hat{\mathbf{M}}'' = \mathbf{0}$ and $\mathbf{P}\hat{\mathbf{\Gamma}}'' = \mathbf{0}$, then the transformations in (A1) will leave them unchanged.

$$\begin{aligned}
v_{gh} &= \text{cov}\left(\frac{\hat{a}}{\hat{c}}, \frac{\hat{b}}{\hat{c}}\right) \\
&\approx E\left[\frac{1}{\hat{c}}\right]^2 \text{cov}(\hat{a}, \hat{b}) + E\left[\frac{1}{\hat{c}}\right] E[\hat{a}] \text{cov}\left(\frac{1}{\hat{c}}, \hat{b}\right) + E\left[\frac{1}{\hat{c}}\right] E[\hat{b}] \text{cov}\left(\hat{a}, \frac{1}{\hat{c}}\right) + E[\hat{a}] E[\hat{b}] \text{cov}\left(\frac{1}{\hat{c}}, \frac{1}{\hat{c}}\right) \\
&\approx E\left[\frac{1}{\hat{c}}\right]^2 \text{cov}(\hat{a}, \hat{b}) - E\left[\frac{1}{\hat{c}}\right] E[\hat{c}]^{-2} E[\hat{a}] \text{cov}(\hat{c}, \hat{b}) - E\left[\frac{1}{\hat{c}}\right] E[\hat{c}]^{-2} E[\hat{b}] \text{cov}(\hat{a}, \hat{c}) \\
&\quad + E[\hat{c}]^{-4} E[\hat{a}] E[\hat{b}] \text{cov}(\hat{c}, \hat{c}) \\
&\approx \left[\frac{\kappa^2}{c^2}\right] \text{cov}(\hat{a}, \hat{b}) - \left[\frac{a}{\kappa c^3}\right] \text{cov}(\hat{c}, \hat{b}) - \left[\frac{b}{\kappa c^3}\right] \text{cov}(\hat{a}, \hat{c}) + \left[\frac{ab}{\kappa^4 c^4}\right] \text{var}(\hat{c}) \\
&\approx \left[\frac{\kappa^2}{c^2}\right] \text{cov}(\hat{a}, \hat{b}) - \left[\frac{a}{\kappa c^2}\right] \text{cov}(\mathbf{P}\hat{\Gamma}'', \hat{b}) - \left[\frac{a}{\kappa c^2}\right] \text{cov}(\hat{a}, \mathbf{P}\hat{\Gamma}'') + \left[\frac{ab}{\kappa^4 c^2}\right] \text{var}(\mathbf{P}\hat{\Gamma}'') \\
&\approx \frac{1}{c^2} \left[\kappa^2 \mathbf{\Pi}_g \mathbf{V}'' \mathbf{\Pi}_h^t - \frac{a}{\kappa} \mathbf{P} \mathbf{\Lambda}''^t \mathbf{\Pi}_h^t - \frac{b}{\kappa} \mathbf{\Pi}_g \mathbf{\Lambda}'' \mathbf{P}^t + \frac{ab}{\kappa^4} \mathbf{P} \mathbf{\Omega}'' \mathbf{P}^t \right]
\end{aligned} \tag{A4}$$

Therefore, the full matrix \mathbf{V}' is:

$$\mathbf{V}' = e^{(-2\mathbf{P}\hat{\Gamma}'')} \left[\kappa^2 \mathbf{\Pi} \mathbf{V}'' \mathbf{\Pi}^t - \kappa^{-1} (\mathbf{\Pi} \mathbf{M}'' \mathbf{P} \mathbf{\Lambda}''^t \mathbf{\Pi}^t + \mathbf{\Pi} \mathbf{\Lambda}'' \mathbf{P}^t \mathbf{M}''^t \mathbf{\Pi}^t) + \kappa^{-4} \mathbf{\Pi} \mathbf{M}'' \mathbf{M}''^t \mathbf{\Pi}^t \mathbf{P} \mathbf{\Omega}'' \mathbf{P}^t \right]. \tag{A5}$$

Similar derivations yield

$$\mathbf{\Lambda}' = e^{(-\mathbf{P}\hat{\Gamma}'')} \left[\kappa \mathbf{\Pi} \mathbf{\Lambda}'' \mathbf{\Pi}^t - \kappa^{-2} \mathbf{\Pi} \mathbf{M}'' \mathbf{P} \mathbf{\Omega}'' \mathbf{P}^t \right] \tag{A6}$$

and

$$\mathbf{\Omega}' = \mathbf{\Pi} \mathbf{\Omega}'' \mathbf{\Pi}^t. \tag{A7}$$

Finally, we obtain estimates of \mathbf{V}' , $\mathbf{\Lambda}'$, and $\mathbf{\Omega}'$ by replacing the terms in (A5), (A6), and (A7) with their sample estimates. Note that $\kappa = 1 + \frac{1}{2} \text{var}(\mathbf{P}\hat{\Gamma}'') = 1 + \frac{1}{2} \mathbf{P} \mathbf{\Omega}'' \mathbf{P}^t \approx 1$ when sample sizes are large (because then $\text{var}(\mathbf{P}\hat{\Gamma}'')$ is very small); in such cases it may be convenient to assume $\kappa = 1$ in Equations (A5) and (A6).

Once we have $\hat{\mathbf{M}}'$, $\hat{\mathbf{\Gamma}}'$, $\hat{\mathbf{V}}'$, $\hat{\mathbf{\Lambda}}'$, and $\hat{\mathbf{\Omega}}'$, we construct $\hat{\mathbf{\Sigma}}'$, the vector of the group standard deviations in the “single prime” metric, as well as the covariance matrices $\hat{\mathbf{Z}}' = \text{Cov}(\hat{\mathbf{M}}', \hat{\mathbf{\Sigma}}')$, and $\hat{\mathbf{W}}' = \text{Cov}(\hat{\mathbf{\Sigma}}', \hat{\mathbf{\Sigma}}')$. The vector of group standard deviations is simply the vector of exponentiated elements of

$\hat{\mathbf{\Gamma}}': \hat{\mathbf{\Sigma}}' = \exp(\hat{\mathbf{\Gamma}}')$. The covariance matrices are then estimated using the Delta method:

$$\begin{aligned}
 \hat{\mathbf{Z}}' &= Cov(\hat{\mathbf{M}}', \exp(\hat{\mathbf{\Gamma}}')) \\
 &\approx Cov(\hat{\mathbf{M}}', \hat{\mathbf{\Gamma}}')[diag(\hat{\mathbf{\Sigma}}')] \\
 &= \hat{\mathbf{\Lambda}}'[diag(\hat{\mathbf{\Sigma}}')],
 \end{aligned} \tag{A8}$$

and

$$\begin{aligned}
 \hat{\mathbf{W}}' &= Cov(\exp(\hat{\mathbf{\Gamma}}'), \exp(\hat{\mathbf{\Gamma}}')) \\
 &\approx [diag(\hat{\mathbf{\Sigma}}')]Cov(\hat{\mathbf{\Gamma}}', \hat{\mathbf{\Gamma}}')[diag(\hat{\mathbf{\Sigma}}')] \\
 &= [diag(\hat{\mathbf{\Sigma}}')]\hat{\mathbf{\Omega}}'[diag(\hat{\mathbf{\Sigma}}')].
 \end{aligned} \tag{A9}$$

Online Appendix B: Standard Errors of Estimated Between-Group Gaps

Equation (10) defines the estimated gap between groups g and h as $\widehat{D}_{gh} = \widehat{d}_{gh}^*/\widehat{s}_{gh}^*$, where \widehat{d}_{gh}^* and \widehat{s}_{gh}^* are, respectively, the difference in the mean value of y^* and the pooled standard deviation of y^* in groups g and h . Let v_{gh}^* and w_{gh}^* be the g, h elements of \mathbf{V}^* and \mathbf{W}^* , respectively. Then the sampling variance of \widehat{d}_{gh}^* will be given by

$$\delta_{gh}^* = \text{var}(\widehat{d}_{gh}^*) = v_{gg}^* + v_{hh}^* - 2v_{gh}^*. \quad (\text{B1})$$

Likewise, the sampling variance of the pooled standard deviation \widehat{s}_{gh}^* will be

$$\begin{aligned} \eta_{gh}^* = \text{var}(\widehat{s}_{gh}^*) &\approx \frac{1}{16s_{gh}^{*2}} (4\sigma_g^{*2}w_{gg}^* + 4\sigma_h^{*2}w_{hh}^* + 8\sigma_g^*\sigma_h^*w_{gh}^*) \\ &= \frac{1}{4s_{gh}^{*2}} (\sigma_g^{*2}w_{gg}^* + \sigma_h^{*2}w_{hh}^* + 2\sigma_g^*\sigma_h^*w_{gh}^*). \end{aligned} \quad (\text{B2})$$

If we assume that the sampling errors in \widehat{d}_{gh}^* and \widehat{s}_{gh}^* are independent (which is not strictly true, since the estimated means and variances on which they depend are jointly estimated, but which will generally be approximately true in modestly large samples), the sampling variance of the gap \widehat{D}_{gh}^* can then be approximated as (Goodman, 1960):

$$\begin{aligned} \text{var}(\widehat{D}_{gh}^*) &= \text{var}\left(\frac{\widehat{d}_{gh}^*}{\widehat{s}_{gh}^*}\right) \\ &\approx \frac{1}{s_{gh}^{*2}} \text{var}(\widehat{d}_{gh}^*) + d_{gh}^{*2} \text{var}\left(\frac{1}{\widehat{s}_{gh}^*}\right) + \text{var}(\widehat{d}_{gh}^*) \cdot \text{var}\left(\frac{1}{\widehat{s}_{gh}^*}\right) \\ &\approx \frac{\delta_{gh}^*}{s_{gh}^{*2}} + \frac{d_{gh}^{*2}\eta_{gh}^*}{s_{gh}^{*4}} + \frac{\delta_{gh}^*\eta_{gh}^*}{s_{gh}^{*4}} \\ &= \frac{\delta_{gh}^*}{s_{gh}^{*2}} \left[1 + D_{gh}^{*2} \frac{\eta_{gh}^*}{\delta_{gh}^*} + \frac{\eta_{gh}^*}{s_{gh}^{*2}} \right]. \end{aligned} \quad (\text{B3})$$

Although it would seem that we can then estimate $var(\widehat{D}_{gh}^*)$ by using the relevant terms from $\widehat{\mathbf{V}}^*$, $\widehat{\mathbf{W}}^*$, $\widehat{\boldsymbol{\Sigma}}^*$, $\widehat{\mathbf{S}}$, and $\widehat{\mathbf{G}}$ in (B3), Goodman (1960) notes that an estimator of $var(\widehat{D}_{gh}^*)$ must account for the fact that the expected values of $\frac{1}{\widehat{s}_{gh}^{*2}}$ and \widehat{d}_{gh}^{*2} will be larger than the desired values in (B3). Taking this into account, a better estimate of $var(\widehat{D}_{gh}^*)$ will be

$$\widehat{var}(\widehat{D}_{gh}^*) \approx \frac{\widehat{\delta}_{gh}^*}{\widehat{s}_{gh}^{*2}} \left[1 + \widehat{D}_{gh}^{*2} \frac{\widehat{\eta}_{gh}^*}{\widehat{\delta}_{gh}^*} - \frac{\widehat{\eta}_{gh}^*}{\widehat{s}_{gh}^{*2}} \right]. \quad (\text{B4})$$

In simulations (not shown), we find that (B4) produces accurate standard errors across the range of conditions in our simulations, with the exception of cases where sample sizes are small ($n = 25$) and the cutscores are poorly located; in those cases, (B4) produced standard errors that were slightly too large, on average.

Online Appendix C: Detailed Simulation Results Tables

Table C1: Bias in Estimated Means, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Low	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		50	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		100	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		400	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
	0.3	25	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		50	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		100	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		400	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
HOMOP	0.0	25	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		50	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		100	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		400	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
	0.3	25	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		50	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		100	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
		400	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C2: RMSE of Estimated Means, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.2502	0.2129	0.2244	0.2061	0.2746	0.2031	0.2161	0.1910
		50	0.1712	0.1492	0.1557	0.1444	0.1783	0.1411	0.1456	0.1328
		100	0.1195	0.1049	0.1092	0.1012	0.1238	0.0995	0.1020	0.0933
		400	0.0592	0.0522	0.0542	0.0502	0.0607	0.0494	0.0506	0.0464
	0.3	25	0.2514	0.2141	0.2267	0.2071	0.2781	0.2035	0.2212	0.1909
		50	0.1710	0.1497	0.1576	0.1445	0.1819	0.1413	0.1489	0.1331
		100	0.1195	0.1054	0.1096	0.1013	0.1241	0.0997	0.1026	0.0935
		400	0.0590	0.0521	0.0540	0.0500	0.0605	0.0493	0.0505	0.0463
HOMOP	0.0	25	0.2214	0.2112	0.2189	0.2034	0.2113	0.1960	0.2051	0.1881
		50	0.1552	0.1482	0.1541	0.1433	0.1454	0.1372	0.1435	0.1317
		100	0.1094	0.1042	0.1087	0.1008	0.1025	0.0970	0.1014	0.0929
		400	0.0545	0.0519	0.0541	0.0501	0.0511	0.0482	0.0506	0.0463
	0.3	25	0.2285	0.2124	0.2194	0.2045	0.2223	0.1977	0.2042	0.1881
		50	0.1644	0.1492	0.1545	0.1435	0.1587	0.1398	0.1436	0.1321
		100	0.1217	0.1057	0.1087	0.1010	0.1204	0.1010	0.1013	0.0934
		400	0.0764	0.0540	0.0540	0.0503	0.0816	0.0558	0.0505	0.0471

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C3: Ratio of Median Estimated Standard Error to Empirical Standard Error for Estimated Means, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.9033	0.9482	1.1378	0.9686	0.9238	0.9381	1.2098	0.9615
		50	0.9518	0.9742	0.9914	0.9809	0.9525	0.9721	0.9979	0.9821
		100	0.9756	0.9879	0.9904	0.9907	0.9700	0.9818	0.9886	0.9901
		400	0.9919	0.9980	0.9971	0.9990	0.9915	0.9956	0.9955	0.9970
	0.3	25	0.9044	0.9467	1.7144	0.9655	0.9609	0.9411	2.3691	0.9632
		50	0.9556	0.9739	1.0797	0.9806	0.9544	0.9737	1.1044	0.9831
		100	0.9774	0.9861	0.9970	0.9899	0.9732	0.9813	1.0075	0.9894
		400	0.9958	1.0030	1.0002	1.0035	0.9989	0.9980	1.0008	0.9995
HOMOP	0.0	25	0.9740	0.9807	0.9855	0.9847	0.9591	0.9762	0.9791	0.9789
		50	0.9854	0.9897	0.9906	0.9888	0.9852	0.9889	0.9914	0.9906
		100	0.9903	0.9955	0.9938	0.9945	0.9902	0.9909	0.9934	0.9944
		400	0.9955	1.0000	0.9979	0.9998	0.9953	0.9975	0.9967	0.9982
	0.3	25	0.9768	0.9889	1.0027	1.0031	0.9618	0.9864	1.0032	1.0030
		50	0.9912	0.9998	1.0084	1.0118	0.9918	1.0017	1.0132	1.0144
		100	0.9994	1.0030	1.0119	1.0161	0.9985	1.0016	1.0164	1.0168
		400	1.0056	1.0147	1.0215	1.0276	1.0065	1.0096	1.0230	1.0231

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C4: Bias in Estimated Standard Deviations, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	-0.0298	-0.0098	-0.0202	0.0012	-0.0423	-0.0132	-0.0297	-0.0038
		50	-0.0116	-0.0037	-0.0035	0.0012	-0.0157	-0.0057	-0.0080	-0.0016
		100	-0.0053	-0.0016	-0.0010	0.0008	-0.0075	-0.0026	-0.0027	-0.0006
		400	-0.0011	-0.0003	-0.0001	0.0003	-0.0016	-0.0005	-0.0005	0.0000
	0.3	25	-0.0332	-0.0122	-0.0333	-0.0002	-0.0454	-0.0142	-0.0385	-0.0041
		50	-0.0129	-0.0049	-0.0112	0.0004	-0.0172	-0.0058	-0.0144	-0.0015
		100	-0.0058	-0.0021	-0.0028	0.0004	-0.0076	-0.0031	-0.0051	-0.0009
		400	-0.0014	-0.0005	-0.0004	0.0001	-0.0016	-0.0006	-0.0006	-0.0001
HOMOP	0.0	25	-0.0071	-0.0046	-0.0064	-0.0029	-0.0118	-0.0077	-0.0097	-0.0058
		50	-0.0032	-0.0021	-0.0030	-0.0013	-0.0053	-0.0040	-0.0049	-0.0030
		100	-0.0016	-0.0010	-0.0015	-0.0006	-0.0027	-0.0019	-0.0024	-0.0014
		400	-0.0003	-0.0002	-0.0003	-0.0001	-0.0005	-0.0004	-0.0005	-0.0003
	0.3	25	0.0014	0.0057	0.0046	0.0080	-0.0079	-0.0001	0.0003	0.0036
		50	0.0054	0.0081	0.0080	0.0095	-0.0006	0.0039	0.0050	0.0065
		100	0.0072	0.0093	0.0096	0.0103	0.0021	0.0053	0.0073	0.0076
		400	0.0083	0.0100	0.0107	0.0107	0.0041	0.0069	0.0091	0.0088

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C5: RMSE of Estimated Standard Deviations, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.2684	0.2218	0.2492	0.1686	0.2722	0.2045	0.2402	0.1542
		50	0.1772	0.1489	0.1373	0.1154	0.1744	0.1379	0.1342	0.1057
		100	0.1216	0.1033	0.0894	0.0799	0.1196	0.0955	0.0843	0.0737
		400	0.0591	0.0507	0.0436	0.0394	0.0585	0.0472	0.0410	0.0364
	0.3	25	0.2727	0.2267	0.2685	0.1706	0.2801	0.2089	0.2495	0.1557
		50	0.1786	0.1524	0.1631	0.1163	0.1778	0.1405	0.1531	0.1062
		100	0.1228	0.1048	0.0979	0.0808	0.1207	0.0978	0.0956	0.0744
		400	0.0602	0.0518	0.0446	0.0399	0.0585	0.0478	0.0419	0.0367
HOMOP	0.0	25	0.0095	0.0074	0.0089	0.0063	0.0168	0.0121	0.0138	0.0106
		50	0.0050	0.0042	0.0049	0.0038	0.0084	0.0073	0.0082	0.0066
		100	0.0030	0.0026	0.0029	0.0024	0.0052	0.0047	0.0052	0.0043
		400	0.0013	0.0012	0.0013	0.0012	0.0023	0.0021	0.0022	0.0020
	0.3	25	0.1477	0.1478	0.1478	0.1479	0.1364	0.1357	0.1358	0.1357
		50	0.1477	0.1478	0.1478	0.1479	0.1356	0.1356	0.1357	0.1357
		100	0.1478	0.1479	0.1479	0.1479	0.1355	0.1356	0.1357	0.1357
		400	0.1478	0.1479	0.1480	0.1480	0.1355	0.1356	0.1357	0.1357

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C6: Ratio of Median Estimated Standard Error to Empirical Standard Error for Estimated Standard Deviations, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.8500	0.8821	0.8741	0.9276	0.8721	0.8882	0.9214	0.9355
		50	0.9236	0.9427	0.9117	0.9605	0.9327	0.9462	0.9028	0.9685
		100	0.9597	0.9703	0.9714	0.9808	0.9647	0.9731	0.9730	0.9848
		400	0.9928	0.9943	0.9933	0.9958	0.9910	0.9926	0.9981	0.9984
	0.3	25	0.8553	0.8883	1.8978	0.9321	0.8812	0.8910	3.8317	0.9369
		50	0.9309	0.9424	0.9177	0.9645	0.9339	0.9463	0.9605	0.9731
		100	0.9643	0.9757	0.9368	0.9826	0.9676	0.9682	0.9471	0.9847
		400	0.9880	0.9921	0.9900	0.9937	0.9972	0.9943	0.9910	0.9981
HOMOP	0.0	25	1.0947	1.0943	1.1013	1.1020	0.8212	0.9730	0.9905	0.9792
		50	1.0893	1.0829	1.0851	1.0671	0.9938	0.9906	0.9974	0.9987
		100	1.0635	1.0403	1.0719	1.0464	0.9940	0.9788	0.9714	0.9990
		400	0.9831	0.9959	0.9790	0.9751	0.9870	0.9761	1.0111	0.9985
	0.3	25	1.0531	1.0919	1.0933	1.0758	0.6942	1.0389	0.9951	1.0083
		50	1.0747	1.0695	1.0763	1.0514	0.9920	1.0286	0.9951	0.9955
		100	1.0533	1.0394	1.0604	1.0370	0.9765	0.9603	0.9844	0.9561
		400	1.0233	1.0132	1.0378	1.0159	1.0145	0.9809	0.9860	0.9870

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C7: Bias in Estimated ICC, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.0245	0.0092	0.0162	0.0085	0.0336	0.0145	0.0276	0.0157
		50	0.0104	0.0044	0.0073	0.0037	0.0137	0.0074	0.0124	0.0078
		100	0.0050	0.0020	0.0035	0.0017	0.0071	0.0036	0.0058	0.0038
		400	0.0010	0.0004	0.0007	0.0002	0.0014	0.0007	0.0013	0.0008
	0.3	25	0.0248	0.0091	0.0149	0.0081	0.0309	0.0127	0.0261	0.0145
		50	0.0104	0.0044	0.0071	0.0037	0.0134	0.0059	0.0122	0.0069
		100	0.0047	0.0020	0.0034	0.0017	0.0062	0.0035	0.0059	0.0037
		400	0.0013	0.0006	0.0009	0.0005	0.0012	0.0007	0.0012	0.0007
HOMOP	0.0	25	0.0142	0.0092	0.0128	0.0061	0.0212	0.0140	0.0175	0.0105
		50	0.0065	0.0043	0.0060	0.0028	0.0096	0.0072	0.0089	0.0054
		100	0.0032	0.0020	0.0029	0.0013	0.0049	0.0034	0.0043	0.0026
		400	0.0006	0.0004	0.0005	0.0002	0.0010	0.0007	0.0010	0.0005
	0.3	25	0.0194	0.0112	0.0132	0.0067	0.0323	0.0188	0.0181	0.0123
		50	0.0116	0.0063	0.0065	0.0036	0.0196	0.0116	0.0097	0.0070
		100	0.0079	0.0038	0.0032	0.0020	0.0147	0.0090	0.0055	0.0049
		400	0.0057	0.0024	0.0011	0.0011	0.0111	0.0061	0.0021	0.0026

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C8: RMSE of Estimated ICC, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.0325	0.0149	0.0205	0.0139	0.0514	0.0228	0.0333	0.0224
		50	0.0139	0.0084	0.0105	0.0079	0.0207	0.0134	0.0171	0.0131
		100	0.0075	0.0052	0.0060	0.0049	0.0122	0.0086	0.0101	0.0082
		400	0.0028	0.0023	0.0026	0.0023	0.0047	0.0038	0.0041	0.0036
	0.3	25	0.0382	0.0152	0.0197	0.0139	0.0495	0.0212	0.0324	0.0214
		50	0.0140	0.0085	0.0105	0.0080	0.0216	0.0125	0.0172	0.0127
		100	0.0074	0.0053	0.0060	0.0049	0.0116	0.0088	0.0101	0.0085
		400	0.0029	0.0024	0.0025	0.0023	0.0049	0.0039	0.0041	0.0037
HOMOP	0.0	25	0.0187	0.0147	0.0176	0.0124	0.0298	0.0216	0.0247	0.0190
		50	0.0099	0.0083	0.0095	0.0074	0.0151	0.0131	0.0146	0.0118
		100	0.0059	0.0052	0.0057	0.0047	0.0094	0.0084	0.0093	0.0077
		400	0.0025	0.0023	0.0025	0.0023	0.0040	0.0038	0.0040	0.0036
	0.3	25	0.0234	0.0160	0.0179	0.0130	0.0406	0.0243	0.0250	0.0197
		50	0.0140	0.0096	0.0100	0.0079	0.0229	0.0157	0.0151	0.0126
		100	0.0094	0.0062	0.0059	0.0050	0.0169	0.0120	0.0098	0.0090
		400	0.0062	0.0034	0.0026	0.0024	0.0117	0.0072	0.0045	0.0045

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Table C9: Ratio of Median Estimated Standard Error to Empirical Standard Error for Estimated ICC, by ICC, CV, Location of Cutscores, Sample Sizes, and Model Type

Model	CV	N	ICC = 0.05				ICC = 0.20			
			Skewed	Mid	Wide	Many	Skewed	Mid	Wide	Many
HETOP	0.0	25	0.9827	1.1873	1.9220	1.1269	0.9055	1.0452	2.9128	1.0017
		50	1.1286	1.1223	1.1092	1.0782	1.0558	1.0288	1.1521	1.0111
		100	1.0934	1.0636	1.0750	1.0535	1.0013	0.9864	0.9774	1.0054
		400	0.9984	1.0024	0.9796	0.9776	0.9946	0.9840	1.0128	1.0013
	0.3	25	0.7358	1.1892	3.4893	1.1227	1.0323	1.1090	4.6791	1.0414
		50	1.1170	1.1238	1.5162	1.0774	1.0330	1.0612	2.1247	1.0085
		100	1.0760	1.0539	1.0869	1.0490	1.0401	0.9720	1.2346	0.9678
		400	1.0140	1.0172	1.0398	1.0176	0.9739	0.9775	0.9899	0.9895
HOMOP	0.0	25	1.0955	1.0947	1.1022	1.1026	0.8303	0.9729	0.9900	0.9788
		50	1.0895	1.0834	1.0857	1.0676	0.9946	0.9910	0.9978	0.9995
		100	1.0635	1.0406	1.0720	1.0466	0.9941	0.9788	0.9717	0.9989
		400	0.9832	0.9959	0.9791	0.9752	0.9871	0.9762	1.0113	0.9987
	0.3	25	1.0537	1.0931	1.0943	1.0766	0.7105	1.0380	0.9960	1.0082
		50	1.0749	1.0701	1.0769	1.0518	0.9916	1.0283	0.9956	0.9954
		100	1.0537	1.0396	1.0609	1.0372	0.9762	0.9602	0.9846	0.9563
		400	1.0235	1.0133	1.0379	1.0160	1.0144	0.9810	0.9862	0.9870

NOTE: CV = coefficient of variation; ICC = intraclass correlation coefficient; Skewed = skewed cutscores 5/30/55; Mid = mid cutscores 20/50/80; Wide = wide cutscores 5/50/95, Many = many cutscores 5/25/50/75/95.

Online Appendix D: Additional Tables for Section 3

Table D1. Descriptive Statistics for NAEP and State Data.

			Group Sample Size					Cut1	Cut2	Cut3	ICC	CV	
	Groups		Min	p25	Median	p75	Max						
NAEP Math	Grade 4	2009	50	1990	2770	2930	3550	7810	0.175	0.602	0.942	0.038	0.139
	Grade 4	2011	50	2710	3210	3590	4520	9880	0.167	0.587	0.933	0.031	0.135
	Grade 8	2009	50	1890	2710	2840	3460	7530	0.263	0.658	0.922	0.042	0.115
	Grade 8	2011	50	2130	2790	2930	3860	8280	0.255	0.645	0.918	0.032	0.105
NAEP Reading	Grade 4	2009	50	2130	2960	3130	3720	8420	0.321	0.666	0.924	0.031	0.169
	Grade 4	2011	50	2750	3290	3720	4490	10140	0.320	0.659	0.922	0.028	0.158
	Grade 8	2009	50	1910	2720	2840	3400	7580	0.236	0.670	0.973	0.031	0.134
	Grade 8	2011	50	2020	2680	2815	3750	7990	0.228	0.658	0.965	0.024	0.119
Average			50	2191	2891	3099	3844	8454	0.246	0.643	0.937	0.032	0.134
State Math	Grade 4	2006	1244	20	48	70	93	345	0.079	0.327	0.795	0.140	0.221
	Grade 5	2006	594	20	68	145	216	432	0.076	0.343	0.810	0.114	0.198
	Grade 6	2006	567	20	76	162	226	563	0.081	0.345	0.785	0.152	0.226
	Grade 7	2006	566	20	77	157	222	518	0.093	0.339	0.785	0.139	0.222
	Grade 8	2006	428	23	104	194	276	668	0.107	0.335	0.781	0.135	0.221
State Reading	Grade 4	2006	1243	20	48	70	93	346	0.032	0.138	0.531	0.117	0.226
	Grade 5	2006	596	20	68	147	218	431	0.012	0.086	0.551	0.102	0.202
	Grade 6	2006	566	20	78	162	229	566	0.024	0.150	0.668	0.115	0.219
	Grade 7	2006	567	20	78	161	226	517	0.019	0.100	0.507	0.108	0.221
	Grade 8	2006	430	20	103	194	280	671	0.011	0.087	0.508	0.101	0.227
Average			680	20	75	146	208	506	0.053	0.225	0.672	0.122	0.218

NOTE: p25 = 25th percentile; p75 = 75th percentile; ICC = intraclass correlation coefficient; CV = coefficient of variation of variances. Sample sizes are rounded to the nearest integer for state data. Sample sizes for NAEP are rounded to the nearest 10 to comply with NCES data reporting requirements.

Table D2. Correlations Between HETOP Estimates and Uncoarsened Score Estimates by Test Subject, Grade and Year.

			Means					Standard Deviations		
			(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
		Estimate 1:	H4	H20	Orig.	H4	H4	H20	Orig.	H4
		Estimate 2:	Orig.	Trans.	Trans.	Trans.	Orig.	Trans.	Trans.	Trans.
NAEP Math	Grade 4	2009	0.998	1.000	1.000	0.998	0.882	0.995	0.966	0.935
	Grade 4	2011	0.997	1.000	0.999	0.998	0.923	0.998	0.963	0.967
	Grade 8	2009	0.997	1.000	1.000	0.997	0.870	0.996	0.973	0.907
	Grade 8	2011	0.997	1.000	1.000	0.997	0.890	0.995	0.978	0.901
NAEP Reading	Grade 4	2009	0.988	1.000	0.999	0.992	0.791	0.998	0.951	0.868
	Grade 4	2011	0.992	1.000	0.998	0.995	0.738	0.994	0.942	0.831
	Grade 8	2009	0.995	1.000	0.999	0.996	0.840	0.992	0.929	0.929
	Grade 8	2011	0.995	1.000	0.999	0.995	0.875	0.994	0.938	0.942
Average			0.995	1.000	0.999	0.996	0.851	0.995	0.955	0.910
State Math	Grade 4	2006	0.988	1.000	0.999	0.988	0.830	0.991	0.974	0.819
	Grade 5	2006	0.990	1.000	1.000	0.990	0.864	0.990	0.990	0.861
	Grade 6	2006	0.990	1.000	0.999	0.991	0.866	0.992	0.968	0.864
	Grade 7	2006	0.991	1.000	1.000	0.992	0.857	0.993	0.983	0.855
	Grade 8	2006	0.986	1.000	1.000	0.986	0.861	0.995	0.989	0.862
State Reading	Grade 4	2006	0.944	0.999	0.998	0.946	0.671	0.980	0.895	0.647
	Grade 5	2006	0.962	1.000	0.999	0.962	0.701	0.986	0.941	0.626
	Grade 6	2006	0.981	1.000	0.999	0.982	0.716	0.984	0.871	0.718
	Grade 7	2006	0.955	1.000	0.998	0.956	0.667	0.982	0.835	0.629
Grade 8	2006	0.941	1.000	0.998	0.941	0.756	0.979	0.872	0.712	
Average			0.973	1.000	0.999	0.973	0.779	0.987	0.932	0.759

NOTE: H4 = heteroskedastic ordered probit model with four proficiency categories based on original 3 cutscores; H20 = heteroskedastic ordered probit model with 20 categories based on 19 equally spaced cutscores. Lowest correlation in each column is indicated in bold.

Online Appendix E: Normality Maximizing Procedure

We would like to assess whether observed test score distributions are respectively normal. Let the observed scale score metric be denoted y . We assume the distributions are respectively normal, then find a monotonic transformation $f(y) = y^*$ that will render the distributions of y^* as near to normal as possible, and then compare the estimated group means and standard deviations in this metric to those estimated from the HETOP model. If they align well, it indicates that the function f has successfully rendered all groups' distributions of y^* normal, which would only be possible under respective normality.

Here we describe a procedure for estimating the function f in order to transform the observed scale scores, and determine whether the means and standard deviations of these transformed scores, denoted $\hat{y}^* = \hat{f}(y)$, more closely match the HETOP estimates than do the means and standard deviations of the original scale scores. In section 3 we applied this procedure separately to each empirical grade-year-subject dataset.

First, we standardized y using the grand mean and standard deviation across all groups.¹⁴ Then, we defined a set of $K = 19$ equally spaced cut points ranging from -2.25 to 2.25 in increments of 0.25. We denote these c_1, c_2, \dots, c_{19} . Note that these cut points are in the original y metric. We coarsened the standardized scale scores into 20 ordered categories using the c_k cut points and treated these as 20 ordered proficiency categories. We fit the HETOP model to the coarsened data, and obtain estimates of $\mathbf{M}^*, \mathbf{\Sigma}^*, \mathbf{C}^*$ (these estimates are referred to as the “HETOP20” estimates in Section 3).

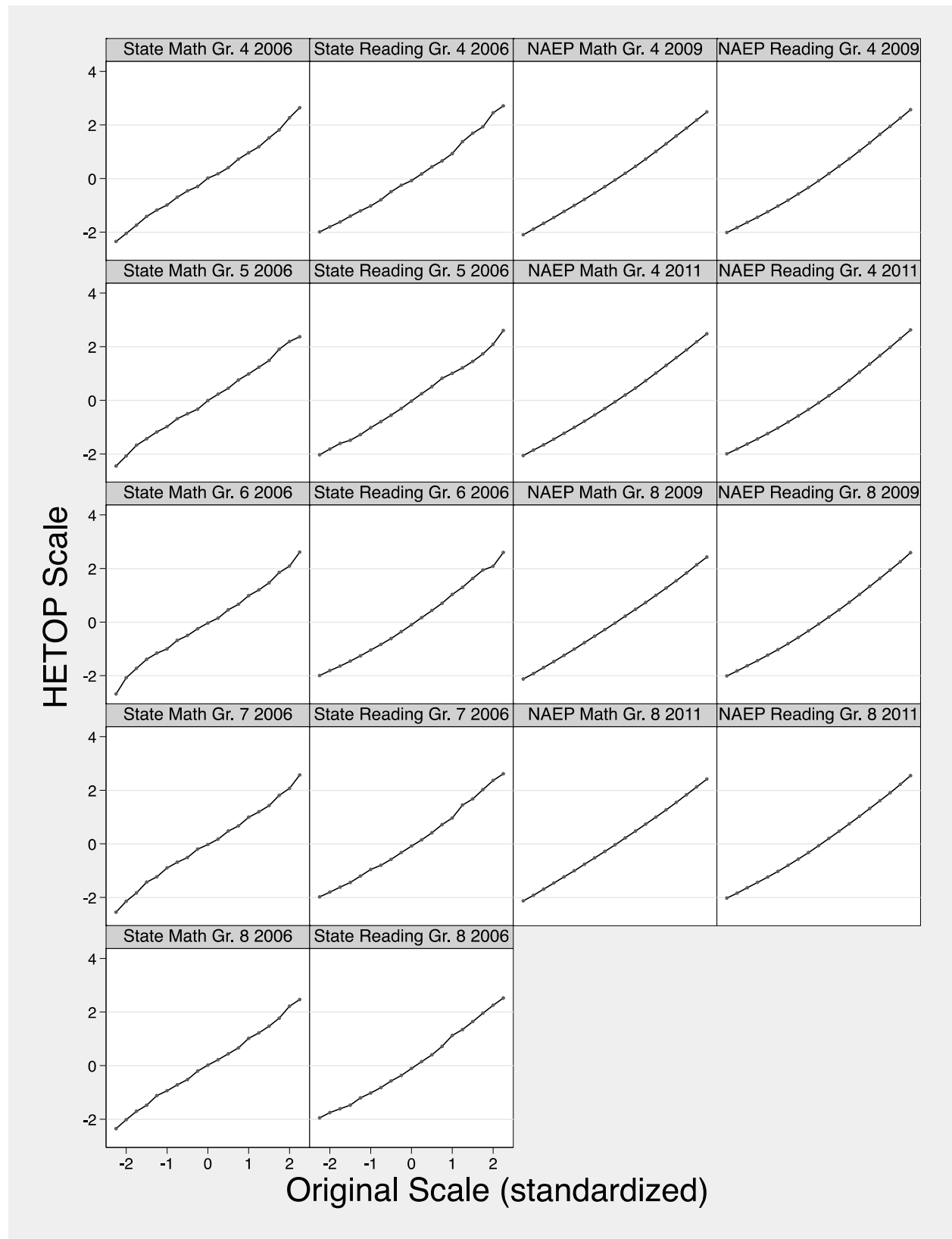
We then plotted values of \hat{c}_k^* against the values c_k and used monotone cubic interpolation (MCI; Fritsch & Carlson, 1980) to estimate a function \hat{f} such that $\hat{f}(c_k) = \hat{c}_k^*$. The use of MCI ensures that \hat{f} is monotonic, differentiable everywhere, and passes through all of the observed pairs (c_k, \hat{c}_k^*) . We then transformed the original scale scores using this MCI transformation, generating values $\hat{y}^* = \hat{f}(y)$. For

¹⁴ We do this only for computational convenience. Because standardizing is a linear transformation applied to all scale scores simultaneously, it does not impact the results of the procedure.

scale scores greater than c_{19} or less than c_1 we used linear extrapolation based on the slope of the estimated transformation at these two endpoints. If the assumption of respective normality is valid, this procedure should render the \hat{y}^* scores normally distributed within each of the G groups. We refer to \hat{f} as a “normality-maximizing” function because even if the distributions are not respectively normal, it attempts to make them simultaneously as nearly normal as possible. Figure E1 displays the original scale score cutscores (after standardization) and the HETOP cutscores for all 18 datasets.

We then computed the means and standard deviations of the transformed scale scores in each group. These are referred to as the as the “transformed” scale score estimates in the normalized metric in Section 3.

Figure E1. Original and Transformed (HETOP) Cutscores for All 18 Datasets.



Online Appendices References

Fritsch, F. N., & Carlson, R. E. (1980). Monotone piecewise cubic interpolation. *SIAM Journal on Numerical Analysis*, 17(2), 238–246. <http://doi.org/10.1137/0717021>

Ho, A. D. (2008). The problem with "proficiency": Limitations of statistics and policy under No Child Left Behind. *Educational Researcher*, 37(6), 351-360.