

Estimating High School GPA Weighting Parameters With a Graded Response Model

John Hansen, *Harvard Graduate School of Education*, Philip Sadler
and Gerhard Sonnert, *Harvard-Smithsonian Center for Astrophysics*

The high school grade point average (GPA) is often adjusted to account for nominal indicators of course rigor, such as “honors” or “advanced placement.” Adjusted GPAs—also known as weighted GPAs—are frequently used for computing students’ rank in class and in the college admission process. Despite the high stakes attached to GPA, weighting policies vary considerably across states and high schools. Previous methods of estimating weighting parameters have used regression models with college course performance as the dependent variable. We discuss and demonstrate the suitability of the graded response model for estimating GPA weighting parameters and evaluating traditional weighting schemes. In our sample, which was limited to self-reported performance in high school mathematics courses, we found that commonly used policies award more than twice the bonus points necessary to create parity for standard and advanced courses.

Keywords: college admission, graded response model, high school grades, item response theory

The high school grade point average (GPA) is a high-stakes measure of academic achievement. In the college admission process, a small change in GPA can have significant consequences, especially at universities using a high school rank or GPA threshold for automatic admission or disqualification (Horn & Flores, 2003). Zimmerman (2014) found that students who barely qualified for admission to public universities in Florida earned more later in life than otherwise similar students who barely missed qualifying based on their GPAs. Shifts toward test-optional and class rank-based admission policies in recent years suggest that the importance of high school grades may be increasing (Hiss & Franks, 2014).

Atkinson and Geiser (2009) concluded that researchers generally agree that academic success in high school coursework is the best predictor of college success. This consensus among scholars is reflected among college admission professionals, who have ranked high school grades in college preparatory courses as the top factor in college admission decisions for decades (Clinedinst, Koranteng, & Nicola, 2016). Yet, there is no consensus about how several years of high school grades across a heterogeneous set of courses should be aggregated into a single number. In California, high school transcripts often include several GPAs, each computed using a different approach: for admission to the University of California or California State University system, advanced

courses receive an extra grade point (University of California, 2016); to determine eligibility for the Cal-Grant, the applicant’s first and last year of high school and all physical education courses are excluded, and no extra points are awarded for advanced courses (Cal-Grant, 2016); to assign class rank within a high school—which has implications for college admission in California—schools are free to choose their own approach. In North Carolina, until 2015, awarding two bonus grade points for advanced placement (AP) courses was a statewide policy (i.e., a B in an AP course received 5.0 grade points, while an A in a standard course received 4.0) (North Carolina State Board of Education, 2015). In Miami-Dade County Schools, AP course grades of A or B receive two bonus points, and a course grade of C receives one bonus point (Miami-Dade County Public Schools, 2016).

Two fundamental problems with using high school GPA as a high-stakes measure of student achievement are (1) not all students take the same courses and (2) grades are ostensibly measured on an ordinal—not interval—scale. As shown by Arrow (1951) and discussed by Vickers (2000) in the context of GPAs, there is no optimal method for aggregating ordinal measures into a single composite. To compare academic achievement across students who took different courses from different teachers at different times, letter grades are typically transformed to numbers and averaged, yielding a GPA. In practice, Lang (2007) found that in practice the most common aggregations were unweighted averages (all courses are weighted equally), weighted averages (e.g., where some courses receive a weight of zero), and rigor-adjusted averages that award extra points for nominal designations of course rigor (e.g., “honors” or “advanced placement”). For a student who took n high school courses, where w_i denotes the scalar weight for course i ’s

John Hansen is a doctoral candidate at Harvard Graduate School of Education, Cambridge, MA 02138; john_hansen@mail.harvard.edu. Philip Sadler is F.W. Wright Senior Lecturer in the Department of Astronomy and Gerhard Sonnert is a Lecturer on Astronomy and Research Associate at Harvard-Smithsonian Center for Astrophysics, Cambridge, MA 02138; psadler@cfa.harvard.edu; gsonnert@cfa.harvard.edu.

contribution to the GPA, one can express the weighted GPA equation as

$$\text{Weighted GPA} = \sum_{i=1}^n w_i (\text{GRADE}_i). \quad (1)$$

The incorporation of a measure of course rigor by awarding bonus points for a subset of courses can be expressed as

$$\text{Rigor Adjusted GPA} = \sum_{i=1}^n w_i (\text{GRADE}_i + \text{BONUS}_i). \quad (2)$$

The term *weighted GPA* is commonly used to refer to rigor-adjusted GPAs, and there is no common designation for a GPA wherein some courses receive weights of zero. Many states and high schools include both a scalar-weight adjustment (e.g., for physical education courses, $w_i = 0$) and a bonus adjustment for advanced courses. For sake of simplicity, hereafter we use the term weighted GPA to refer collectively to GPAs with course-varying scalar weights and/or bonus point adjustments.

In this study, we discuss and demonstrate the use of item response theory (IRT) for examining GPA weighting policies. We used the graded response model (Samejima, 1969) because its underlying assumptions are compatible with conventional GPAs, it uses the same data as conventional GPAs, and graded response model (GRM) parameters are conceptually similar to conventional GPA weighting parameters. This article is organized into four sections. Section 1 briefly reviews previous research on the validity of the high school GPA. Section 2 discusses previous examples of IRT applications to course grade data and the theoretical suitability of using a GRM to study high school GPA weighting parameters. In Section 3, we fit a GRM to actual high school course grade data and estimate bonus point parameters for advanced courses on a conventional GPA scale. Section 4 concludes.

Previous Research on the Validity of the High School GPA

Many studies of the high school grade point average have shown that high school GPA predicts first-year college grades (Kobrin, Patterson, Shaw, Mattern, & Barbuti, 2008), cumulative college GPA (Geiser & Santelices, 2007), and rates of college graduation (Bowen, Chingos, & McPherson, 2009). Studies have also found that students who take advanced courses in high school, such as AP courses, tend to outperform students who do not (Keng & Dodd, 2008). Geiser and Santelices (2004), Klopfenstein and Thomas (2005), Sadler and Tai (2007b), Warne (2017), and Warne, Nagaishi, Slade, Hermesmeyer, and Peck (2014) found that the strength of the relationship between advanced course participation and college achievement depends on the additional covariates in the model. An unresolved question is the correct methodological approach to estimating how GPAs should be weighted by course type.

Criterion-Based Validation of Weights for Advanced Courses

Previous researchers have relied on an external criterion to estimate GPA weights for advanced courses. Although they reached different conclusions, Geiser and Santelices (2004) and Sadler and Tai (2007a, b) employed similar, criterion-based empirical strategies. They used multiple regression

where the outcome was a measure of college achievement, and the independent variables included indicators for participation in advanced coursework. Sadler and Tai (2007b) based their recommendations of .50 bonus grade points for honors courses and 1.0 bonus grade points for AP courses on the relationship between high school science course participation and college science course performance. They compared the final grades in an introductory college science course among students who did and those who did not take honors or AP versions of the same course in high school. Using college chemistry as an example, they fit the following ordinary least squares (OLS) regression model:

$$\begin{aligned} \text{COL_CHEM} = & \beta_0 + \beta_1(\text{HS_GRADE}) + \beta_2(\text{HON_CHEM}) \\ & + \beta_3(\text{AP_CHEM}) + \epsilon, \end{aligned} \quad (3)$$

where *COL_CHEM* is the grade in introductory college chemistry, *HON_CHEM* and *AP_CHEM* are dichotomous indicator variables, and standard high school chemistry is the comparison category. *HS_GRADE* is the letter grade earned in a student's last high school chemistry course. In some specifications, they also added indicators for college instructor, demographic control variables, and SAT scores. The college course performance differentials uniquely explained by honors and AP course performance were translated to a grade scale by identifying the magnitude of a high school grade (*HS_GRADE*) increase equivalent to the differential. Using coefficients from a regression represented by Equation 3, the recommended bonus point adjustment for taking an honors chemistry course was

$$\text{Bonus} = \frac{\beta_{\text{HON_CHEM}}}{\beta_{\text{HS_GRADE}}}. \quad (4)$$

A fundamental challenge for the Sadler and Tai (2007a, b) methodology is that one must select an external criterion. For example, one could use all college grades, first-year college grades, total credits accumulated, persistence beyond the first year, or degree completion. In selecting college science course performance as a criterion, Sadler and Tai (2007a, b) implicitly assume that college science course performance suffices as a scale for measuring the equivalency between $\beta_{\text{HON_CHEM}}$ and $\beta_{\text{HS_GRADE}}$. In this framework, the correct high school grade adjustment is the one that minimizes prediction error in college course performance, either because one cares about predicting college course performance for its own sake or because college course performance is considered a sufficient measure of a broader construct of interest. Either way, the assumption is that high school grades are meaningful insofar as they predict college grades in the courses students choose to take, at the colleges where they are admitted and choose to enroll.

An additional challenge for criterion selection is choosing what kind of variability in the criterion is the "correct" variability for identifying grade weights. Because many factors other than advanced course participation can explain variability in college performance, researchers have tried to isolate variance attributable to advanced course participation by using covariates to restrict the criterion variance of interest to within-college variance (Klopfenstein & Thomas, 2005), within-college-by-instructor variance (Sadler & Tai, 2007b), and/or variance unexplained by demographics (Rothstein,

2004). Rothstein (2004) discussed selection bias and modeling strategies in detail, especially the issue of nonrandom selection of students into colleges and courses.

Ultimately, criterion and model specification are consequential because many factors are correlated with high school grades, participation in advanced courses, and college outcomes. Different model specifications support different inferences and different interpretations of the high school GPA construct. These challenges motivate the use of item response theory (IRT), because item response models rely only on observed measures of high school performance—the exact data one needs to compute a conventional GPA—and nothing more.

Item Response Model GPAs

Statistical methods for adjusting grades received in different courses predate the development of item response theory. Linn (1966) and Young (1993) reviewed methods proposed in the twentieth century, which were almost exclusively non-IRT approaches. Young (1990a) published the first study that used an item response model to link grades from different courses onto a common scale. Studying college performance among undergraduates, he found that the correlation between college GPA and SAT scores increased substantially when using a graded response model-scaled college “GPA” instead of the conventional college GPA. Similar to Young (1990a), recent studies of course-adjusted GPA methods have generally focused on using item response theory to correct the college grade point average or class rank for differential grading stringency across college courses and departments (Bailey, Rosenthal, & Yoon, 2016; Caulkins, Larkey, & Wei, 1996; Johnson, 1997; Lei, Bassiri, & Schultz, 2001; Young, 1990b). An exception is Bassiri and Schultz (2003), who used a Rasch rating scale model to create a universal scale of high school difficulty by using ACT performance as an anchor.

The present study’s application of item response theory is distinct from previous work because it focuses explicitly on course-level parameters. In this respect, the present study is perhaps most similar to Korobko, Glas, Bosker, and Luyten (2008), who used IRT to compare the difficulty of university entrance exams in The Netherlands, where students select different subsets of exams. Our study illustrates item response theory’s use for constructing a scale of course difficulty, which could be used to evaluate policies about awarding bonus grade points for advanced course participation.

Unlike most previous studies, we did not explicitly seek an estimator of student skill that optimally predicted academic success in college. We chose not to incorporate standardized achievement measures (e.g., SAT or AP exam scores) in our models, because such tests plausibly measure a construct that differs from the one measured by the high school GPA (Willingham, Pollack, & Lewis, 2002). If certain standardized test scores indeed measured the same construct as the high school GPA (or the construct that one believes the high school GPA ought to measure), including the scores in the model could improve estimation of grade-by-course difficulty parameters and improve the comparability of grades received at different schools. For example, the simultaneous inclusion of school-level intercept parameters and students’ AP exam scores could potentially identify between-school differences in grading stringency. Schools where students received higher AP exam scores than one would expect given their course grades

would be identified as grading more stringently. This model would allow one to measure student academic skill in a way that simultaneously takes into account course grades, school grading stringency, and AP exam performance—assuming that AP exam scores measured the same construct as course grades, and assuming that GPAs should be adjusted for between-school differences. Ultimately, the validity of GPAs estimated from such models depends on how one defines the GPA construct and the potential consequences of their use (Kane, 2013). If one defines GPA as an absolute measure of curricular knowledge, then AP exams could be beneficial. If one defines GPA as a measure of academic performance relative to local standards, then AP exams could be distortional. We omitted standardized test scores to emphasize that additional data are neither necessary nor necessarily desirable for estimating course difficulty parameters. College admission policies regularly treat grades received from different high schools as interchangeable, and omitting standardized test scores allowed us to estimate course difficulty on a common scale under the assumption of interchangeable courses.

The appeal of estimating high school course difficulty with item response theory is that one can use the exact same data as conventional GPAs and no additional data, retain the key assumptions of conventional GPAs (1–3 below), and relax nonessential constraints (4–6 below). Like conventional GPAs, IRT can measure high school achievement in a way that (1) treats the construct of interest underlying all of a student’s grades (hereafter θ) as one-dimensional; (2) assumes that students who receive higher grades in a given course tend to have higher levels of θ than students who receive lower grades in the same course; and (3) assumes that the construct of interest is identifiable from observed grades received in courses taken, and that course-taking patterns—including nonexistent or “missing” grade data—can be ignored. In item response theory, this is the conditional independence assumption. Rubin (1976) classifies data that are missing at random conditional on observable data as missing at random (MAR). Assuming data are MAR implies that a student who has no grade in a certain course does not necessarily have a lower or higher θ than one who has a grade in the course. Setting aside their potential flaws, we treat these assumptions as fundamental assumptions of GPAs. In short, the trait of interest is unidimensional, higher grades imply higher trait values, and the trait is identifiable from observed data only. A potential benefit of estimating a “GPA” with item response theory is that one can relax additional constraints imposed by conventional GPAs. In the case of the GRM, (4) differences in course difficulty are estimated, not presupposed; (5) the trait’s scale is estimated from the data without assuming an interval property (e.g., the difference on the scale between an A and a B is not assumed to be the same as the difference between a B and a C, and these distances are not assumed to be constant across courses); and (6) all courses are not assumed to be equally informative measures of θ . Theoretically, relaxing these constraints will improve the internal coherence of the measure unless the imposed constraints indeed conform to the data-generating process.

To be clear, the GRM’s statistical properties make it a good candidate for examining the properties of conventional approaches to calculating GPAs and addressing policy-relevant questions, such as how advanced courses should be weighted. However, similar to other IRT models, the GRM’s technical complexity and latent scale make it a less appealing

candidate for widespread use as an alternative GPA-weighting algorithm. The GRM estimates parameters with conventional GPA weighting analogs, but it does not directly estimate them on the conventional GPA scale of Equations 1 or 2. The GRM assumes the existence of a latent variable, conventionally assumed to be normally distributed with a mean of zero and unit variance. In estimating parameters, the model finds the θ scale using a cumulative log-odds principle (Ostini & Nering, 2005). Formally, the equation for estimating the GRM can be expressed as

$$P(Y_{ij} \geq k | a_i, b_{ik}, \theta_j) = \frac{1}{1 + \exp\{-a_i(\theta_j - b_{ik})\}}, \quad (5)$$

where Y_{ij} is the letter grade earned by student j in course i ; K represents the possible letter grades; a_i can be interpreted as the precision with which course i measures θ ; θ is a latent variable one could interpret as student academic skill (or the construct one supposes is measured by grades in high school courses); and b_{ik} is a measure of the difficulty of achieving letter grade k in course i . For each course, $K - 1$ threshold parameters are estimated. Each b_{ik} is the θ value in which a student for whom $\theta_j = b_{ik}$ has a .50 fitted probability of earning letter grade k or higher in course i . Comparing estimates of b_{ik} allows one to compare difficulties of earning letter grades across different courses. True equivalence in difficulty also requires letter grades across courses to have identical values of a_i . In this case, fitted probabilities of exceeding the thresholds would be equivalent for all values of θ , not only where $\theta_j = b_{ik}$. More precisely, a_i parameterizes the rate at which, for course i , the probability of exceeding thresholds

changes as a function of θ . If the probability of exceeding thresholds changes little as a function of θ for course i , the lower value of a_i indicates that student grades in course i are less informative measures of θ .

GRM Parameter Interpretation in the Weighted GPA Context

Figure 1 is a stylized representation of perfect correspondence between a common weighted GPA policy and grade-by-course difficulty parameters from a graded response model. The most common policy for adjusting grades in honors and AP courses is to award an additional letter grade for AP courses and one half of a letter grade for honors courses (Lang, 2007), which we plot in panel A. Panel B of Figure 1 illustrates how grade-by-course difficulty parameters can be converted between a weighted GPA scale and the GRM's latent variable scale. For sake of illustration, we assume a θ scale where the average student earns a B or higher in 50% of standard courses. In this scenario, where the GRM estimates conform perfectly to the weighted GPA policy, the GRM boundary parameters (b_{ik}) for common mathematics and science courses are a linear shift of -3.0 units from the weighted GPA scale. Under these stylized conditions, note that one can transform θ to a 5.0 GPA scale simply by adding three points and truncating GPAs above 5.0 or below zero.

Panel B illustrates the relationship between boundary parameters (b_{ik}) and grade-by-course difficulty. In this example, the b parameters for a B in AP Statistics and an A in standard statistics are the same. Because $b_B^{AP_STATS} = b_A^{STATS}$, the curves intersect where $p = .50$. For student j such that $\theta_j = b_B^{AP_STATS} = b_A^{STATS}$, the probability of earning a grade of

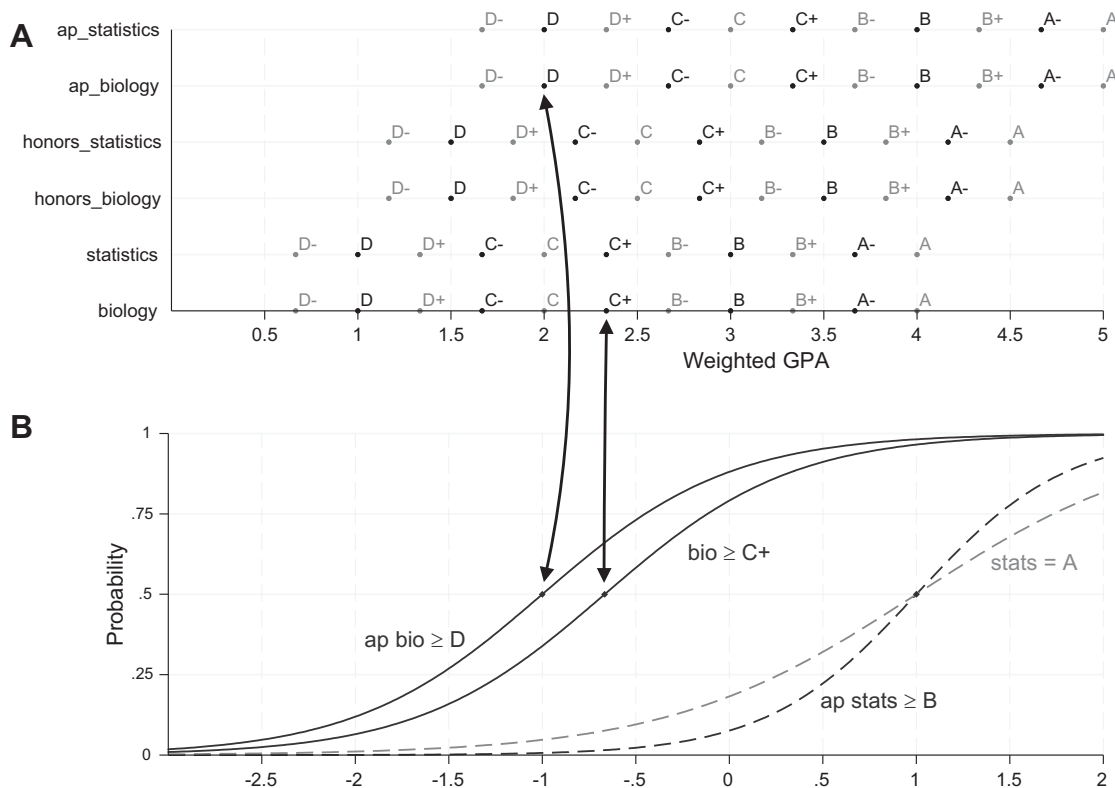


FIGURE 1. Conceptual model for GRM-estimation of high school GPA weighting parameters at the letter-grade-by-course level. (A) Weighted GPA policy in which honors courses receive .50 bonus grade points and advanced placement courses receive 1.0 bonus grade points. (B) Probability of earning a given course grade or higher as a function of θ .

B or higher in AP Statistics is the same as the probability of earning an A in a standard statistics course ($p = .50$). Hence, these grade-by-course combinations are equally difficult to achieve for student j .

Panel B also illustrates that the a_i parameters are different for standard statistics and AP Statistics courses, which means that whether it is more or less difficult to earn a B in AP Statistics or an A in standard statistics depends on θ . Nevertheless, the differential is symmetric around b , implying equivalent difficulty on average (or, more precisely, equivalent average difficulty for a symmetric distribution of students centered at b). Biology illustrates a case where the a_i parameters are the same for AP Biology and standard biology but the b parameters differ by .33 units of θ . As a result, for any fixed value of θ , the probability of earning a D or higher in AP Biology is greater than the probability of earning a C+ in standard biology.

In the example above, a difference of .33 units on the θ scale is equivalent to .33 grade points on a 5-point GPA scale. In general, the magnitude of similar differences on the θ -scale will be sensitive to rescaling, even though the scale itself is arbitrary up to a linear transformation. As shown by Lord (1980), consider a transformation of the current θ scale to an alternative scale with greater variance, such that the new scale has m times the standard deviation of the old ($m\sigma_{OLD} = \sigma_{NEW}$). Fitted probabilities are preserved so long as a_i and b_{ik} are rescaled as well.

$$P(Y_{ij} \geq k | a_i, b_{ik}, \theta_j) = \frac{1}{1 + \exp\{-a_i(\theta_j - b_{ik})\}} \\ = \frac{1}{1 + \exp\{\frac{-a_i}{m}(m(\theta_j - b_{ik}))\}}. \quad (6)$$

Graphically, this is equivalent to multiplying by m the x -axis labels in Figure 2. Yet, despite the substantive equivalence of the models, the magnitude of the difference on the θ -scale is

not preserved: $(mb_D^{AP-Bio} - mb_{C+}^{Bio}) \neq (b_D^{AP-Bio} - b_{C+}^{Bio})$. We revisit this issue in Section 3.2.

As always, substantive interpretations of parameter estimates rely on the assumption that the GRM model fits the data. As Thissen (2016) argues, the answer to the “bad question” of whether an item response model fits is, in binary terms, that it does not. Certainly, in this case, there are many reasons why Y_{ij} might not be exclusively a function of item parameters and θ . A GRM does not ameliorate all shortcomings of grades or conventional GPA aggregation approaches. Many potential criticisms of the GRM approach apply equally to conventional GPAs. Unless explicitly modeled, neither accounts for variability in grading standards in nominally identical courses offered by different teachers, in different high schools, or on different occasions. The extent to which these omissions are flaws or features depends on the relevant inference and what one believes student GPA is supposed to measure. Concerns that apply similarly to GRM and conventional GPAs—and possible ameliorations—merit additional study, but we do not address them here. Our key point is that a GRM approach is well suited to estimating how various grade-by-course combinations can be compared to one another on a common scale.

Overall, using a GRM to estimate letter grade-by-course difficulty parameters on a latent scale can, at an abstract level, be characterized as a linking study with nonequivalent groups, a common-item design, and concurrent calibration (Kolen & Brennan, 2004). Simultaneous estimation of course parameters (concurrent calibration) requires the existence of multiple courses and crossing of students and courses. Courses with the same name are treated as common items, allowing all courses and students to be linked to a common scale. While courses taught at different schools are obviously nonequivalent in many respects, weighted GPA policies regularly treat these courses equivalently. Consequently, differences in course difficulty associated with

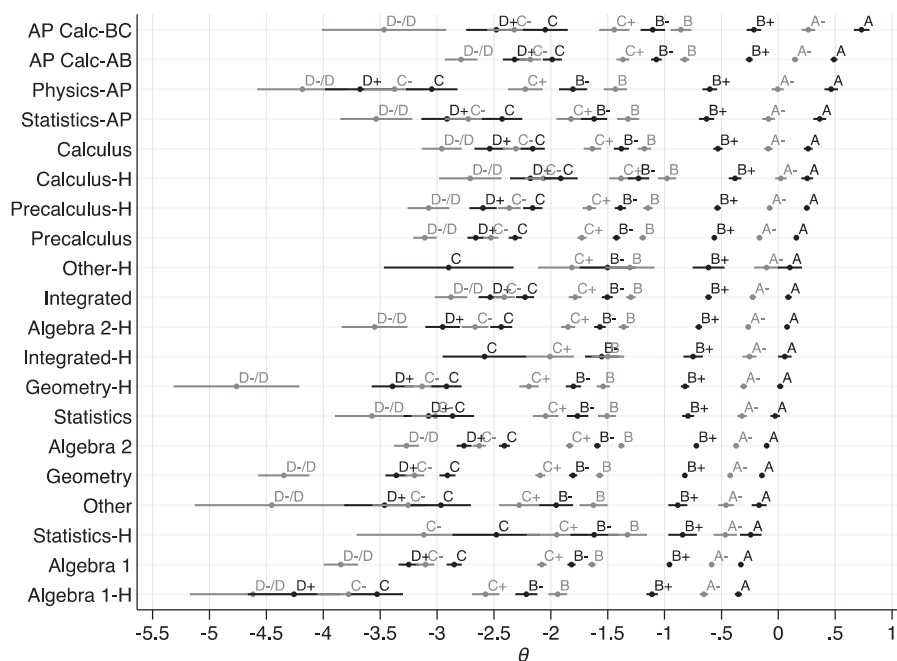


FIGURE 2. Each point indicates the mathematics skill (θ) necessary to have a .50 probability of earning the given letter grade or higher in the given course. Confidence intervals are ± 1 SE. Grades of D- and D are collapsed into a single category due to sparse data at the low end of the grade distribution (see Table A1 in the online Supporting Information).

nominal course identifiers are the policy-relevant quantities of interest. Interpretation of b_{ik} parameters where high school courses are the “items” is analogous to canonical GRM applications under a few conditions, which are discussed in the online Supporting Information.

Demonstration

In this section, we fit a GRM to high school course grade data and use the results to estimate bonus point (i.e., weighting) parameters for advanced courses on a conventional GPA scale.

Data and the GRM Model

The data for this section of the study are from the Factors Influencing College Success in Mathematics survey (FICS-Math), which was gathered in a very similar fashion as the data used by Sadler and Tai (2007a, b). The FICSMath survey was administered to 10,492 first-year college calculus students from 135 randomly selected U.S. colleges during the 2009–2010 academic year. The sample was intended to be approximately—but not precisely—representative of first-year college calculus students. The sample included public, private, large, small, selective, nonselective, religious, and nonreligious institutions from across the United States. Twenty-eight institutions had “Community College” in their name, and six institutions were classified as “Most Competitive” in *Barron’s Profiles of American Colleges* (Barron’s 2009).

Students took the FICSMath survey at the beginning of the course. At the end of the term, the calculus instructor recorded each student’s grade from the course. Most of the survey items focused on pedagogical decisions made by teachers that could contribute to college success in mathematics and science. It also included items about the student’s high school, the student’s demographic background, and, most importantly for this study, the mathematics courses they took in high school and the letter grades they received in those courses. We focused only on quantitative courses (mathematics courses and AP Physics), because the survey has limited data on nonquantitative coursework and our model assumes a unidimensional construct.

We fit the graded response model using the “irt grm” command in Stata 14 (StataCorp, 2015). The default max-

imum likelihood estimator approximates integrals with mean and variance adaptive Gauss–Hermite quadrature. Initial attempts to fit the model encountered some convergence challenges. In response, we collapsed D– and D grades into a single category, because D– grades were relatively rare in our sample (A1). At this point, the algorithm converged within two minutes without issuing errors or warnings. The 853 students who did not report any course grades were excluded from estimation. Next, we estimated the posterior mean of θ for each student using empirical Bayes. This is the GRM analogue to GPA.

Next, to create a sample for subsequent analyses, in this order, we listwise deleted students for whom no final grade was reported (359), students who reported taking fewer than three quantitative courses (1,535), and students who did not respond to the survey item on parental educational attainment (262). The goal of these restrictions was to create a consistent sample that included an external criterion, an approximate measure of high school GPA, and a measure of student socioeconomic status (SES). Table 1 shows that the students taking calculus as first-year college students were relatively high-performing high school students. The average unweighted GPA for quantitative courses in the analytical sample was 3.50.

Table 1 also shows correlations among variables. Similar to previous studies, the IRT estimate of student academic skill— $\hat{\theta}$, which we also refer to as a GRM-weighted GPA—had a stronger correlation with the relevant criterion outcome (.352) than the unweighted GPA (.326). The conventionally weighted GPAs (1.0 bonus for AP courses, .50 bonus for honors courses) correlation with college calculus grade was .347, nearly as strong as $\hat{\theta}$. At the same time, parental educational attainment, a measure of student SES, was positively correlated with calculus grade ($r = .085$) and more strongly correlated with the conventionally weighted GPA ($r = .086$) than the GRM-weighted GPA or unweighted GPA ($r = .042$). This covariance pattern illustrates a concern with using correlations to validate a particular weighting scheme. The weighted GPA may outperform the unweighted GPA in part due to more educated parents encouraging their children to take more advanced courses, and additional support that more educated students provide their children during college. Table A2 shows that the IRT estimate of student academic skill also outperformed the other GPAs in multiple regression models

Table 1. Descriptive Statistics and Correlations for the Analytical Sample ($N = 8,336$)

	Correlations								Mean	SD
	Calc.	GRM GPA	GPA	W. GPA	Tot.	Hon.	AP	Par. BA		
Calculus grade									79.8	14.4
GPA: GRM - $\hat{\theta}$.352								.03	.88
GPA: no weights	.326	.937							3.50	.50
GPA: weighted	.347	.882	.913						3.73	.58
Total course grades	.157	.210	.174	.279					4.77	1.25
Hon. course grades	.130	.176	.144	.469	.279				1.43	1.77
AP course grades	.145	.151	.088	.341	.596	.259			.49	.80
Parent has bach.	.085	.042	.049	.086	.085	.102	.070		.63	.63
SAT-math*	.270	.339	.308	.387	.263	.237	.244	.172	612.8	98.1

Note: All data above except calculus grade were self-reported by students participating in the Factors Influencing College Success in Mathematics study (FICSMath). First-year college calculus classes were sampled from a randomly selected set of U.S. colleges during the 2009–2010 academic year. The weighted GPA awarded 1.0 bonus points for AP courses, and .50 points for honors courses. All correlations were statistically significant ($\alpha < .001$). *SAT-math includes ACT Math scores translated to an SAT scale using an SAT/ACT concordance. SAT and ACT scores were missing for 1,315 students.

that accounted for student SAT score, parental educational attainment, and college course instructor. Essentially, stronger correlations with a criterion and weaker correlations with potential confounds (e.g., parental education) help establish credibility for the model's estimates of grade-by-course difficulty, but, for reasons previously discussed, we do not view these correlations as the definitive arbiters of GPA weighting schemes. The relatively small difference in the correlations with college calculus grade belie substantial differences between the conventional weighted GPA's bonus parameters for advanced courses and the IRT-based bonus parameter estimates.

Figure 2 plots the GRM estimates of grade-by-course difficulty (b_{ik}) on the θ scale. For this sample, which was restricted to students enrolling in first-year college calculus, the ranking of difficulty parameter estimates generally supported the practice of GPA weights for AP courses—particularly for AP Calculus. AP Calculus AB and BC both cover foundational college-level calculus concepts, but BC includes more advanced material. The lower rate at which students with high grades in other classes received As in AP Calculus is the GRM's empirical basis for identifying the course as particularly difficult. The proportion of As received was lowest in AP Calculus BC and second lowest in AP Calculus AB (Table A1 in the online Supporting Information). Students who reported taking either AP Calculus course had higher average grades in their other courses ($M = 3.65, SD = .40$) than students who did not take AP Calculus ($M = 3.47, SD = .52, t(8,334) = 15.14, p < .001$).

Overall, the results presented in Figure 2 support the practice of GPA adjustments of some magnitude for nominal indicators of course difficulty. However, in contrast to weighting schemes in which AP courses receive one full letter grade adjustment, earning a B+ or better in the most difficult AP course appears similarly difficult as earning an A in a standard algebra course (for students with θ around $-.25$). In most cases, it appears that approximately one-third of a letter grade is closer to the course letter grade adjustment to

create equivalent letter grade difficulties between AP and non-AP courses. For honors courses, the typical differential appears to be less than one third of a letter grade.

Linking Latent-Scale Parameters to a GPA Scale

The $BONUS_i$ parameters of interest from Equation 2 are not directly estimated by the GRM, but the GRM's θ -scale estimates of grade-by-course difficulty, \hat{b}_{ik}^θ , can be transformed to a conventional GPA scale and then used to estimate $BONUS_i$ parameters. We used a linear linking (Kolen & Brennan, 2004) approach to identify the linear transformation that mapped the \hat{b}_{ik}^θ for standard courses to a conventional 4.0 GPA, and applied the same linear transformation to honors and advanced courses (additional technical details in the online Supporting Information). Figure 3 displays the results of the transformation: a linear shift of the θ -scale estimates from Figure 2 onto a GPA scale. We only used \hat{b}_{ik}^θ for grades in the A–C range, which meant that \hat{b}_{ik}^{GPA} for standard courses were centered at 3.0. We excluded \hat{b}_{ik}^θ for grades below a C primarily to simplify the demonstration, because \hat{b}_{ik}^θ for grades below a C were estimated imprecisely, and in some cases could not be estimated at all. Substantively, one could also argue that higher-achieving students take more AP courses, and it therefore makes more sense to prioritize accurate grade links in the range of grades that students tend to receive. Ultimately, a detailed treatment of various technical and substantive issues implied by link choices is beyond the scope of this article.

The difference between an unweighted grade point value (4.0 grade points for an A in AP Statistics) and the estimated grade point value (4.25 grade points for an A in AP Statistics) is our estimate for $BONUS_{AP,STATS}$. The largest estimated bonus point adjustment was .45 grade points for AP Calculus BC. One could use the results in Figure 3 to design a weighting system where each course was weighted differently, or one could use the average $BONUS_i$ for AP courses

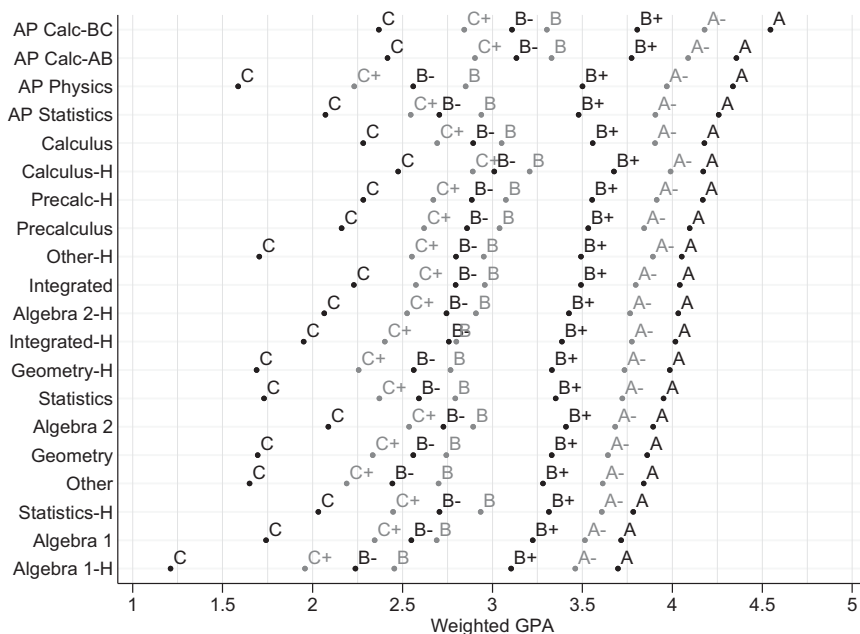


FIGURE 3. Point estimates of grade by course difficulty on a weighted high school GPA scale. Grades below C were estimated imprecisely due to sparse data, and were omitted from the linking process (see the online Supporting Information for details).

to design a weighting system where all AP courses receive the same adjustment, $BONUS_{AP}$. We estimated that the average GPA-scale difference in difficulty between standard and AP courses, $BONUS_{AP}$, was .25 letter grade points—a bit shy of the difference between an A and an A–. Our estimate of $BONUS_{HON}$ was .02 letter grade points—virtually no difference.

Limitations

The purpose of this study is to discuss and demonstrate the GRM approach to evaluating policies for weighting high school GPAs, not to make a definitive claim about the “correct” high school GPA weighting parameters across all U.S. schools. Our study included only one IRT method for estimating latent scale parameters, only one method for transforming latent scale parameters to GPA scale parameters, and only one method for aggregating $BONUS_i$ into $BONUS_{AP}$. Our decisions prioritized clarity in illustration of the methodology over other factors. The effect of these methodological decisions on parameter estimates merits further study.

Another limitation is our sample, which relied on self-reported high school grades in quantitative coursework for students who chose to take college calculus. How accurate are self-reports compared to administrative data sets? Rosen, Porter, and Rogers (2017) compared self-reported math course enrollment and grades to high school students’ transcripts. They characterized the overall level of misreporting as “manageable” (p. 12), and found that inaccurate reporting was mainly an issue among lower-performing students, such as students who received low grades in Algebra I. College calculus enrollees comprised our sample, which suggests that inaccurate reporting may not be substantial. As for focusing only on quantitative coursework, this is arguably an asset in terms of satisfying the model assumption of a unidimensional construct, but a limitation in terms of generalizability, because the results may be different for advanced courses in nonquantitative subjects. A final limitation is that our sample only contained students who chose to take college calculus, a relatively high-achieving sample. An attractive feature of estimating high school GPA parameters with an item response model is that, theoretically, parameter estimates are sample-invariant. If the model fits the data, parameters estimated with high-achieving and low-achieving samples will be linear transformations of one another. Conducting similar analyses with other samples is an important next step for research on high school GPA weighting policy. The National Center for Education Statistics makes data publicly available for samples of students that differ from ours in several important ways: students were sampled within high schools; subgroups of students who may be underrepresented in college calculus courses were proportionately or overrepresented; and actual transcripts, rather than self-reports, were used to identify courses and grades.

Concluding Remarks

Overall, high school grades play a powerful role in education. They motivate students to study, provide feedback to students about their academic performance, and inform college admission committees about students’ high school performance. Despite inconsistencies in grading practices across courses, teachers, and schools, grades tend to predict college success as well as—if not better than—standardized test

scores do. As a result, in recent years, many colleges have placed greater emphasis on high school grades in the college admission process.

To account for differences in grading standards between standard and advanced courses, GPAs are often adjusted to account for nominal indicators of course rigor, such as “advanced placement.” This study discussed and demonstrated how item response theory could be used to estimate course difficulty differentials between courses. The ability to compare the difficulty of various grade-by-course combinations on a conventional weighted GPA scale can support well-informed policy decisions for weighting high school GPAs. If the sole purpose of awarding bonus grade points for advanced course participation is to equate grades from standard and advanced high school courses, our results indicate that widely used policies award excess points for advanced courses. Whether our results generalize beyond our data set and their sensitivity to plausible methodological alternatives is an important question for future research.

Acknowledgments

This research was supported by Grant No. 0813702 from the National Science Foundation. Any opinions, findings, and conclusions in this article are the authors’ and do not necessarily reflect the views of the National Science Foundation. Without the excellent contributions of many people, the FICSMATH project would not have been possible. This article benefitted from feedback from Professor Andrew Ho and students in his course on educational measurement. The authors thank two anonymous reviewers, whose comments and suggestions significantly improved the article.

References

- Arrow, K. J. (1951). *Social choice and individual welfare*. New York, NY: John Wiley and Sons.
- Atkinson, R. C., & Geiser, S. (2009). Reflections on a century of college admissions tests. *Educational Researcher*, 38, 665–676.
- Bailey, M. A., Rosenthal, J. S., & Yoon, A. H. (2016). Grades and incentives: Assessing competing grade point average measures and postgraduate outcomes. *Studies in Higher Education*, 41, 1548–1562.
- Barron’s. (2009). *Barron’s profiles of American colleges 2009*. Hap-pauge, NY: Barron’s.
- Bassiri, D., & Schulz, E. M. (2003). Constructing a universal scale of high school course difficulty. *Journal of Educational Measurement*, 40, 147–161.
- Bowen, W. G., Chingos, M. M., & McPherson, M. S. (2009). *Crossing the finish line: Completing college at America’s public universities*. Princeton, NJ: Princeton University Press.
- Cal-Grant. (2016, Sept. 21). Worksheet to calculate Cal-Grant GPA. Retrieved from http://www.oakparkusd.org/cms/lib5/CA01000794/Centricity/Domain/165/Cal-Grant_GPA_Calculation.pdf
- Caulkins, J. P., Larkey, P. D., & Wei, J. (1996). Adjusting GPA to reflect course difficulty. Retrieved from <http://repository.cmu.edu/cgi/viewcontent.cgi?article=1042&context=heinzworks>
- Clinedinst, M., Koranteng, A.-M., & Nicola, T. (2016). *State of college admission*. Arlington, VA: National Association for College Admission Counseling. Retrieved from <https://indd.adobe.com/view/c555ca95-5bef-44f6-9a9b-6325942ff7cb>
- Geiser, S., & Santelices, M. V. (2004). The role of advanced placement and honors courses in college admissions. Research & Occasional Paper Series CSHE.4.04. Berkeley: University of California Center for Studies in Higher Education. Retrieved from <http://www.cshe.berkeley.edu/sites/default/files/shared/publications/docs/ROP.Geiser.4.04.pdf>

- Geiser, S., & Santelices, M. V. (2007). Validity of high-school grades in predicting student success beyond the freshman year. Research & Occasional Paper Series CSHE.6.07. Berkeley: University of California Center for Studies in Higher Education. Retrieved from <http://files.eric.ed.gov/fulltext/ED502858.pdf>
- Hiss, W. C., & Franks, V. W. (2014). *Defining promise: Optional standardized testing policies in American college and university admissions*. Report of the National Association for College Admission Counseling (NACAC). Retrieved from http://www.nacacnet.org/research/research-data/nacac-research/Documents/Defining_Promise.pdf
- Horn, C. L., & Flores, S. M. (2003). *Executive summary: Percent plans in college admissions. A comparative analysis of three states' experiences*. Retrieved from <https://civilrightsproject.ucla.edu/research/college-access/admissions/percent-plans-in-college-admissions-a-comparative-analysis-of-three-states2019-experiences/>
- Johnson, V.E. (1997). An alternative to traditional GPA for evaluating student performance. *Statistical Science*, 12, 215–269.
- Kane, M. T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50, 1–73.
- Keng, L., & Dodd, B. G. (2008). *A comparison of college performances of AP and non-AP student groups in 10 subject areas*. New York, NY: College Board.
- Klopfenstein, K., & Thomas, M. K. (2005, January). *The advanced placement performance advantage: Fact or fiction?* Paper presented at the meeting of the American Economic Association, Philadelphia, PA. Retrieved from https://www.aeaweb.org/assa/2005/0108_1015_0302.pdf
- Kobrin, J. L., Patterson, B. F., Shaw, E. J., Mattern, K. D., & Barbuti, S. M. (2008). *Validity of the SAT for predicting first-year college grade point average*. Research Report No. 2008-5. New York, NY: College Board.
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating, scaling, and linking*. New York, NY: Springer.
- Korobko, O. B., Glas, C. A., Bosker, R. J., & Luyten, J. W. (2008). Comparing the difficulty of examination subjects with item response theory. *Journal of Educational Measurement*, 45, 139–157.
- Lang, D. M. (2007). Class rank, GPA, and valedictorians: How high schools rank students. *American Secondary Education*, 35(2), 36–48.
- Lei, P. W., Bassiri, D., & Schultz, E. M. (2001). *Alternatives to the grade point average as a measure of academic achievement in college*. Report No. TM033669. Iowa City, IA: American College Testing Program.
- Linn, R. L. (1966). Grade adjustments for prediction of academic performance: A review. *Journal of Educational Measurement*, 3, 313–329.
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Hillsdale, NJ: Lawrence Erlbaum.
- Miami-Dade County Public Schools. (2016) *Rank in class—Grade point average*. School Board rule 6Gx13-5B-1.061. Retrieved from <http://www.dadeschools.net/schoolboard/rules/Chapt5/5b-1.061.pdf>
- North Carolina State Board of Education. (2015). Policy outlining standards to be incorporated into the electronically generated high school transcript. GS 116-11(10a). Retrieved from <http://sbepolicy.dpi.state.nc.us/policies/GCS-L-004.asp?pri=01&cat=L&pol=004&acr=GCS>
- Ostini, R., & Nering, M. L. (2005). *Polytomous item response theory models*. Thousand Oaks, CA: Sage.
- Rosen, J. A., Porter, S. R., & Rogers, J. (2017). Understanding student self-reports of academic performance and course-taking behavior. *AERA Open*, 3(2), 1–14.
- Rothstein, J. M. (2004). College performance predictions and the SAT. *Journal of Econometrics*, 121(1), 297–317.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63, 581–592.
- Sadler, P. M., & Tai, R. H. (2007a). Accounting for advanced high school coursework in college admission decisions. *College and University*, 82(4), 7–14.
- Sadler, P. M., & Tai, R. H. (2007b). Weighting for recognition: Accounting for advanced placement and honors courses when calculating high school grade point average. *NASSP Bulletin*, 91(1), 5–32.
- Samejima, F. (1969). Estimation of latent ability using a response pattern of graded scores. *Psychometrika Monographs*, 24(2, No. 17).
- StataCorp (2015). *Stata statistical software: Release 14*. College Station, TX: StataCorp LP.
- Thissen, D. (2016). Bad questions: An essay involving item response theory. *Journal of Educational and Behavioral Statistics*, 41(1), 81–89.
- University of California. (2016, Sept. 21). Calculating the UC/CSU GPA. Retrieved from http://collegetools.berkeley.edu/documents/cat_113-128/Calculating_GPA.pdf
- Vickers, J. M. (2000). Justice and truth in grades and their averages. *Research in Higher Education*, 41(2), 141–164.
- Warne, R. T. (2017). Research on the academic benefits of the Advanced Placement program: Taking stock and looking forward. *SAGE Open*, 7(1), 1–17.
- Warne, R. T., Nagaishi, C., Slade, M. K., Hermesmeier, P., & Peck, E. K. (2014). Comparing weighted and unweighted grade point averages in predicting college success of diverse and low-income college students. *NASSP Bulletin*, 98(4), 261–279.
- Willingham, W. W., Pollack, J. M., & Lewis, C. (2002). Grades and test scores: Accounting for observed differences. *Journal of Educational Measurement*, 39, 1–37.
- Young, J. W. (1990a). Adjusting the cumulative GPA using item response theory. *Journal of Educational Measurement*, 27, 175–186.
- Young, J. W. (1990b). Are validity coefficients understated due to correctable defects in the GPA? *Research in Higher Education*, 31(4), 319–325.
- Young, J. W. (1993). Grade adjustment methods. *Review of Educational Research*, 63(2), 151–165.
- Zimmerman, S. D. (2014). The returns to college admission for academically marginal students. *Journal of Labor Economics*, 32, 711–754.

Supporting Information

Additional Supporting Information, including statistical code, may be found in the online version of this article at the publisher's website:

Table A1. Letter Grades Received by Course

Table A2. OLS Regression Prediction of Final Grade in College Calculus by GPA

Figure A1. Comparison across scales of estimated b_{ik} parameters in calculus and statistics for standard and Advanced Placement courses. For calculus, AP Calculus AB is plotted. (A) θ scale. (B) GPA scale.

Figure A2. Comparison of b_{ik} parameters estimated from a graded response model in which a parameters are allowed to vary by course (y axis) and a graded response model in which a parameters are constrained to be equal for all courses (x axis).