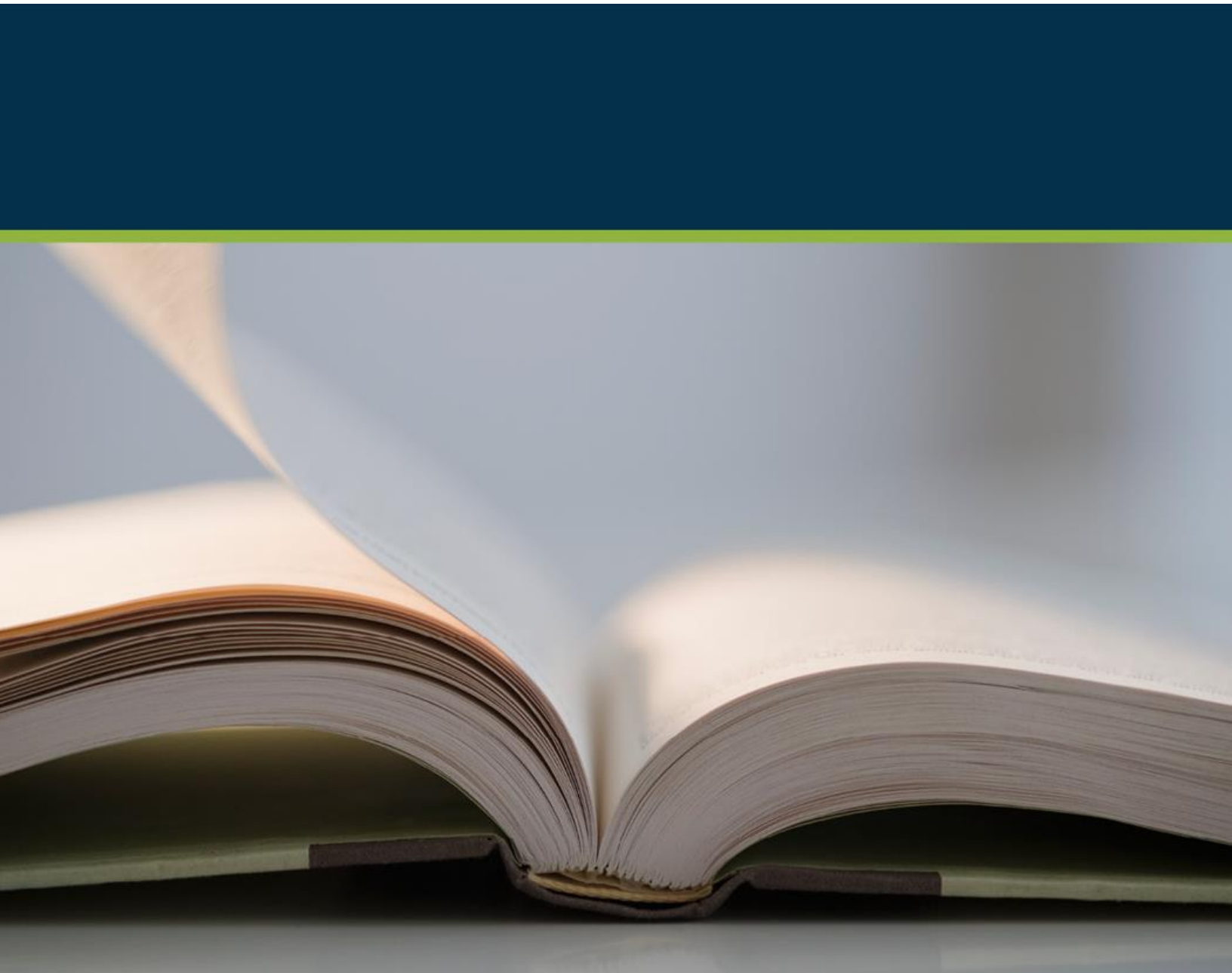


SAS® EVAAS

Statistical Models and Business Rules

Prepared for the North Carolina Department of Public Instruction



Contents

- 1 Introduction to North Carolina’s Value-Added Reporting1**
- 2 Statistical Models2**
 - 2.1 Overview of Statistical Models2
 - 2.2 Gain Model3
 - 2.2.1 Overview3
 - 2.2.2 Why the Gain Model is Needed4
 - 2.2.3 Common Scale in the Gain Model6
 - 2.2.4 Technical Description of the Gain Model9
 - 2.3 Predictive Model.....17
 - 2.3.1 Overview17
 - 2.3.2 Conceptual Explanation17
 - 2.3.3 Technical Description of the District, School, and Teacher Models20
 - 2.4 Projection Model22
 - 2.4.1 Overview22
 - 2.4.2 Technical Description.....23
 - 2.5 Outputs from the Models24
 - 2.5.1 Gain Model.....24
 - 2.5.2 Predictive Model.....25
 - 2.5.3 Projection Model26
- 3 Expected Growth28**
 - 3.1 Overview28
 - 3.2 Technical Description.....28
 - 3.3 Illustrated Example28
- 4 Classifying Growth into Categories31**
 - 4.1 Overview31
 - 4.2 Use Standard Errors Derived from the Models31
 - 4.3 Define Growth Indicators in Terms of Standard Errors31
 - 4.4 Illustrated Examples of Categories31
 - 4.5 Rounding and Truncating Rules33
- 5 Composite Growth Measures34**
 - 5.1 Teacher Composites34
 - 5.1.1 Overview34
 - 5.1.2 Technical Description of the Composite Index Based on Gain Model Measures.....35
 - 5.1.3 Technical Description of the Composite Index Based on Predictive Model Measures.....37
 - 5.1.4 Technical Description of the Combined Composite Index Across Subjects Based on the Gain and Predictive Models.....37
 - 5.2 District and School Composites38
- 6 Input Data Used in the North Carolina Growth Model39**
 - 6.1 Assessment Data Used in North Carolina39
 - 6.2 Student Information39
 - 6.3 Teacher Information40

7 Business Rules	41
7.1 Assessment Verification for Use in Growth Models.....	41
7.1.1 Stretch.....	41
7.1.2 Relevance.....	41
7.1.3 Reliability.....	41
7.2 Pre-Analytic Processing	42
7.2.1 Missing Grade	42
7.2.2 Duplicate (Same) Scores	42
7.2.3 Students with Missing Districts or Schools for Some Scores but Not Others	42
7.2.4 Students with Multiple (Different) Scores in the Same Testing Administration	42
7.2.5 Students with Multiple Grade Levels in the Same Subject in the Same Year	42
7.2.6 Students with Records That Have Unexpected Grade Level Changes	42
7.2.7 Students with Historical EOC Records that Occur Too Early or Late Given the Student’s Testing History	43
7.2.8 Students with Records at Multiple Schools in the Same Test Period.....	43
7.2.9 Outliers.....	43
7.2.10 Linking Records Over Time	44
7.3 Growth Models.....	45
7.3.1 Students Included in the Analysis	45
7.3.2 Minimum Number of Students to Receive a Report	45
7.3.3 Student-Teacher Linkages.....	47

1 Introduction to North Carolina’s Value-Added Reporting

The term “value-added” refers to a statistical analysis used to measure students’ academic growth. Conceptually and as a simple explanation, value-added or growth measures are calculated by comparing the exiting achievement to the entering achievement for a group of students. Although the concept of growth is easy to understand, the implementation of a growth model is more complex.

First, there is not just one growth model; there are multiple growth models depending on the assessment, students included in the analysis, and level of reporting (district, school, or teacher). For each of these models, there are business rules to ensure the growth measures reflect the policies and practices selected by the State of North Carolina.

Second, in order to provide reliable growth measures, growth models must overcome non-trivial complexities of working with student assessment data. For example, students do not have the same entering achievement, students do not have the same set of prior test scores, and all assessments have measurement error because they are estimates of student knowledge. EVAAS growth models have been in use and available to educators in states since the early 1990s. These growth models were among the first in the nation to use sophisticated statistical models that addressed these concerns.

Third, the growth measures are relative to students’ expected growth, which is in turn determined by the growth that is observed within the actual population of North Carolina test-takers in a subject, grade, and year. Interpreting the growth measures in terms of their distance from expected growth provides a more nuanced, and statistically robust, interpretation.

With these complexities in mind, the purpose of this document is to guide you through North Carolina’s value-added modeling based on the statistical models, business rules, policies, and practices selected by the State of North Carolina and currently implemented by EVAAS. This document includes details and decisions in the following areas:

- Conceptual and technical explanations of analytic models
- Definition of expected growth
- Classifying growth into categories for interpretation
- Explanation of district, school, and teacher composites
- Input data
- Business rules

The state of North Carolina has provided EVAAS growth measures to North Carolina districts, schools, and teachers since 2005. By 2006, district and school value-added reporting was available statewide, and in 2008, Teacher Value-Added reports also became available for parts of the state. The first year of statewide implementation for teacher value-added reporting that included all teachers with students taking the state assessments in grades 4–8 was 2011.

These reports are delivered through the EVAAS web application available at <http://ncdpi.sas.com>. Although the underlying statistical models and business rules supporting these reports are sophisticated and comprehensive, the web reports are designed to be user-friendly and visual so that educators and administrators can quickly identify strengths and opportunities for improvement and then use these insights to inform curricular, instructional, and planning supports.

2 Statistical Models

2.1 Overview of Statistical Models

The conceptual explanation of value-added reporting is simple: compare students' exiting achievement with their entering achievement over two points in time. In practice, however, measuring student growth is more complex. Students start the school year at different levels of achievement. Some students move around and have missing test scores. Students might have "good" test days or "bad" test days. Tests, standards, and scales change over time. A simple comparison of test scores from one year to the next does not incorporate these complexities. However, a more robust value-added model, such as the one used in North Carolina, can account for these complexities and scenarios.

North Carolina's value-added models offer the following advantages:

- **The models use multiple subjects and years of data.** This approach minimizes the influence of measurement error inherent in all academic assessments.
- **The models can accommodate students with missing test scores.** This approach means that more students are included in the model and represented in the growth measures. Furthermore, because certain students are more likely to have missing test scores than others, this approach provides less biased growth measures than growth models that cannot accommodate student with missing test scores.
- **The models can accommodate tests on different scales.** This approach gives flexibility to policymakers to change assessments as needed without a disruption in reporting. It permits more tests to receive growth measures, particularly those that are not tested every year.
- **The models can accommodate team teaching or other shared instructional practices.** This approach provides a more accurate and precise reflection of student learning among classrooms.

These advantages provide robust and reliable growth measures to districts, schools, and teachers. This means that the models provide valid estimates of growth given the common challenges of testing data. The models also provide measures of precision along with the individual growth estimates taking into account all of this information.

Furthermore, because this robust modeling approach uses multiple years of test scores for each student and includes students who are missing test scores, EVAAS value-added measures typically have very low correlations with student characteristics. It is not necessary to make *direct* adjustments for student socioeconomic status or demographic flags because each student serves as their own control. In other words, to the extent that background influences persist over time, these influences are already represented in the student's data. As a 2004 study by The Education Trust stated, specifically with regard to the EVAAS modeling:

[I]f a student's family background, aptitude, motivation, or any other possible factor has resulted in low achievement and minimal learning growth in the past, all that is taken into account when the system calculates the teacher's contribution to student growth in the present.

Source: Carey, Kevin. 2004. "The Real Value of Teachers: Using New Information about Teacher Effectiveness to Close the Achievement Gap." *Thinking K-16* 8(1):27.

In other words, although technically feasible, adjusting for student characteristics in sophisticated modeling approaches is typically not necessary from a statistical perspective, and the value-added reporting in North Carolina does not make any direct adjustments for students' socioeconomic or demographic characteristics. Through this approach, the North Carolina Department of Public Instruction does not provide growth models to educators based on differential expectations for groups of students based on their backgrounds.

Based on North Carolina's state assessment program, there are two approaches to providing growth measures.

- **The gain model (also known as the multivariate response model or MRM)** is used for tests given in consecutive grades, such as EOG Mathematics and Reading assessments in grades 3–7 and the Early Literacy assessments in grades K-2.
- **The predictive model (also known as univariate response model or URM)** is used when a test is given in non-consecutive grades or when performance from previous tests is used to predict performance on another test. This includes EOG Mathematics and Reading assessments in grade 8, EOG Science assessments in grades 5 and 8, end-of-course (EOC) exams, Career and Technical Education (CTE) tests, and college readiness assessments such as PSAT, SAT, and ACT.

There is another model, which is similar to the predictive model except that it is intended as an instructional tool for educators serving students who have not yet taken an assessment.

- **The projection model** is used for all assessments and provides a probability of obtaining a particular score or higher on a given assessment for individual students.

The following sections provide technical explanations of the models. The online Help within the EVAAS web application is available at <https://ncdpi.sas.com>, and it provides educator-focused descriptions of the models.

2.2 Gain Model

2.2.1 Overview

The gain model measures growth between two points in time for a group of students; this is the case for tests given in consecutive grades such as EOG Mathematics and Reading assessments in grades 3–7 and the Early Literacy in grades K-2.¹ **More specifically, the gain model measures the change in relative achievement for a group of students based on the statewide achievement from one point in time to the next.** For state summative assessments, growth is typically measured from one year to the next using the available consecutive grade assessments. For Early Literacy assessments, growth is measured from the beginning of the year to the end of the year within the same grade for first and second grades and from the middle of the year to the end of the year within the same grade for kindergarten. Expected growth means that students maintained their relative achievement among the population of test-takers, and more details are available in [Section 3](#).

¹ Starting with the 2022-23 reporting, growth measures for EOG Mathematics and Reading in grade 8 no longer use the gain model. However, the description here applies to previous years' reporting for those assessments.

There are three separate analyses for EVAAS reporting based on the gain model: one each for districts, schools, and teachers. The district and school models are essentially the same; they perform well with the large numbers of students characteristic of districts and most schools. The teacher model uses a version adapted to the smaller numbers of students typically found in teachers' classrooms.

In statistical terms, the gain model is known as a linear mixed model and can be further described as a multivariate repeated measures model. These models have been used for value-added analysis for almost three decades, but their use in other industries goes back much further. These models were developed to use in fields with very large longitudinal data sets that tend to have missing data.

Value-added experts consider the gain model to be among one of the most statistically robust and reliable models. The references below include foundational studies by experts from RAND Corporation, a non-profit research organization:

- On the **choice of a complex value-added model**: McCaffrey, Daniel F., and J.R. Lockwood. 2008. "Value-Added Models: Analytic Issues." Prepared for the National Research Council and the National Academy of Education, Board on Testing and Accountability Workshop on Value-Added Modeling, Nov. 13-14, 2008, Washington, DC.
- On the **advantages of the longitudinal, mixed model approach**: Lockwood, J.R. and Daniel McCaffrey. 2007. "Controlling for Individual Heterogeneity in Longitudinal Models, with Applications to Student Achievement." *Electronic Journal of Statistics* 1:223-252.
- On the **insufficiency of simple value-added models**: McCaffrey, Daniel F., B. Han, and J.R. Lockwood. 2008. "From Data to Bonuses: A Case Study of the Issues Related to Awarding Teachers Pay on the Basis of the Students' Progress." Presented at Performance Incentives: Their Growing Impact on American K-12 Education, Feb. 28-29, 2008, National Center on Performance Incentives at Vanderbilt University.

2.2.2 Why the Gain Model is Needed

A common question is why growth cannot be measured with a simple gain model that measures the difference between the current year's scores and prior year's scores for a group of students. The example in Figure 1 illustrates why a simple approach is problematic.

Assume that 10 students are given a test in two different years with the results shown in Figure 1. The goal is to measure academic growth (gain) from one year to the next. Two simple approaches are to calculate the mean of the differences *or* to calculate the differences of the means. When there is no missing data, these two simple methods provide the same answer (5.8 on the left in Figure 1). When there is missing data, each method provides a different result (6.9 versus 4.6 on the right in Figure 1).

Figure 1: Scores without Missing Data, and Scores with Missing Data

Student	Previous Score	Current Score	Gain
1	51.9	74.8	22.9
2	37.9	46.5	8.6
3	55.9	61.3	5.4
4	52.7	47.0	-5.7
5	53.6	50.4	-3.2
6	23.0	35.9	12.9
7	78.6	77.8	-0.8
8	61.2	64.7	3.5
9	47.3	40.6	-6.7
10	37.8	58.9	21.1
Column Mean	50.0	55.8	5.8
Difference between Current and Previous Score Means			5.8

Student	Previous Score	Current Score	Gain
1	51.9	74.8	22.9
2		46.5	
3	55.9	61.3	5.4
4		47.0	
5	53.6	50.4	-3.2
6	23.0	35.9	12.9
7	78.6	77.8	-0.8
8	61.2	64.7	3.5
9	47.3	40.6	-6.7
10	37.8	58.9	21.1
Column Mean	51.2	55.8	6.9
Difference between Current and Previous Score Means			4.6

A more sophisticated model can account for the missing data and provide a more reliable estimate of the gain. As a brief overview, the gain model uses the correlation between current and previous scores in the non-missing data to estimate means for all previous and current scores as if there were no missing data. It does this without explicitly imputing values for the missing scores. The difference between these two estimated means is an estimate of the average gain for this group of students. In this example, the gain model calculates the estimated difference to be 5.8. Even in a small example such as this, the estimated difference is much closer to the difference with no missing data than either measure obtained by the mean of the differences (6.9) or the difference of the means (4.6). This method of estimation has been shown, on average, to outperform both of the simple methods.² This small example only considered two grades and one subject for 10 students. Larger data sets, such as those used in the actual value-added analyses for the state, provide better correlation estimates by having more student data, subjects, and grades. In turn, these provide better estimates of means and gains.

This simple example illustrates the need for a model that will accommodate incomplete data sets, which all student testing sets are. The next few sections provide more technical details about how the gain model calculates student growth.

¹ See, for example, S. Paul Wright, "Advantages of a Multivariate Longitudinal Approach to Educational Value-Added Assessment without Imputation," Paper presented at National Evaluation Institute, 2004. Available online at <https://evaas.sas.com/support/EVAAS-AdvantagesOfAMultivariateLongitudinalApproach.pdf>.

2.2.3 Common Scale in the Gain Model

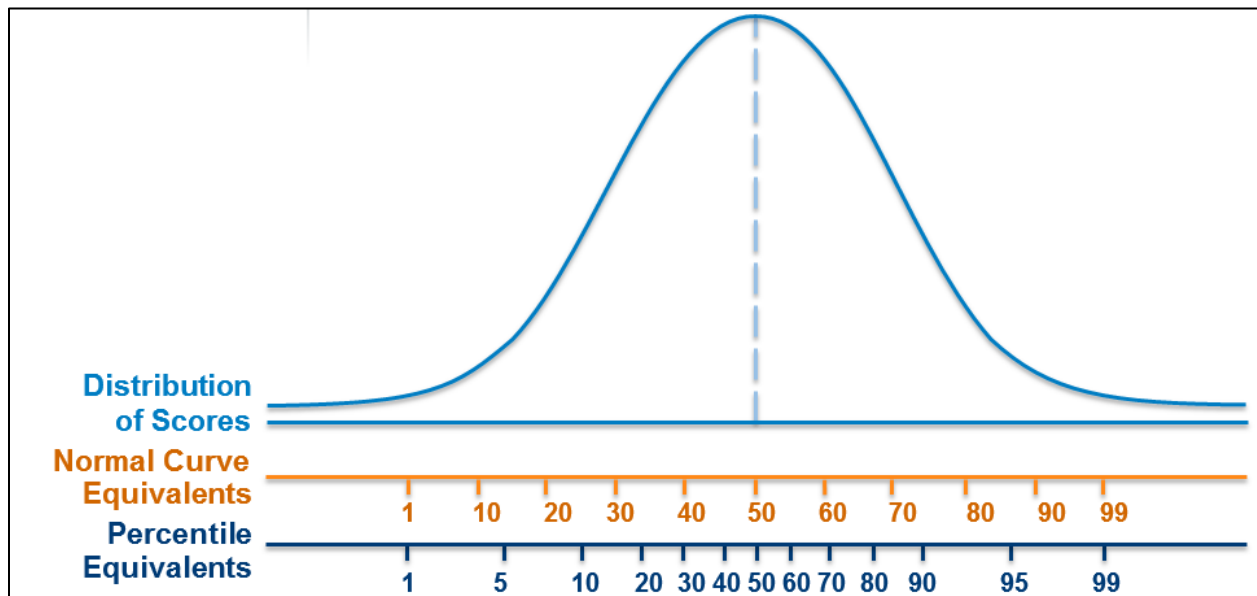
2.2.3.1 Why the Model Uses Normal Curve Equivalents

The gain model estimates academic growth as a “gain,” or the difference between two measures of achievement from one point in time to the next. For such a difference to be meaningful, the two measures of achievement (that is, the two tests whose means are being estimated) must measure academic achievement on a common scale. Even for some vertically scaled tests, there can be different growth expectations for students based on their entering achievement. A reliable alternative whether tests are vertically scaled is to convert scale scores to normal curve equivalents (NCEs).

An NCE distribution is similar to a percentile one. Both distributions provide context as to whether a score is relatively high or low compared to the other scores in the distribution. In fact, NCEs are constructed to be equivalent to percentile ranks at 1, 50 and 99 and to have a mean of 50 and standard deviation of approximately 21.063.

However, NCEs have a critical advantage over percentiles for measuring growth: NCEs are on an equal-interval scale. This means that for NCEs, unlike percentile ranks, the distance between 50 and 60 is the same as the distance between 80 and 90. This difference between the distributions is evident below in Figure 2.

Figure 2: Distribution of Achievement: Scores, NCEs and Percentile Rankings



Furthermore, percentile ranks are usually truncated below 1 and above 99, and NCEs can range below 0 and above 100 to preserve the equal-interval property of the distribution and to avoid truncating the test scale. In a typical year among North Carolina’s state assessments, the average maximum NCE is typically around 115. Although the gain model does not use truncated values, which could create an artificial floor or ceiling in students’ test scores, the web reporting might show NCEs as integers from 1 to 99 for display purposes.

Each NCE distribution is based on a specific assessment, test, subject, and time point. For example, the NCE distribution for 2023 EOG Math in grade 5 is constructed separately from the NCE distribution for 2023 EOG Math in grade 4. The Early Literacy assessments for K-2 have their own NCE distribution as well.

2.2.3.2 Sample Scenario: How to Calculate NCEs in the Gain Model

The NCE distributions used in the gain model are based on a reference distribution of test scores in North Carolina. This reference distribution is the distribution of scores on a state-mandated test for all students in a given year. By definition, the mean (or average) NCE score for the reference distribution is 50 for each grade and subject. For identifying the other NCEs, the gain model uses a method that does not assume that the underlying scale is normal. This method ensures an equal-interval scale, even if the testing scales are not normally distributed.

Table 1 provides an example of how the gain model converts scale scores to NCEs. The first five columns of the table are based on a tabulated distribution of about 115,000 test scores from North Carolina data. In a given subject, grade, and year, the tabulation shows, for each given score, the number of students who scored that score (“Frequency”) as well as the percentage (“Percent”) that frequency represents out of the entire population of test-takers. The table also tabulates the “Cumulative Frequency as the number of students who made that score or lower and its associated percentage (“Cumulative Percent”).

The next column, “Percentile Rank,” converts each score to a percentile rank. As a sample calculation using the data in Table 1 below, the score of 425 has a percentile rank of 45.2. The data show that 43.5% of students scored *below* 425 while 46.9% of students scored *at or below* 425. To calculate percentile ranks with discrete data, the usual convention is to consider half of the 3.4% reported in the Percent column to be “below” the cumulative percent and “half” above the cumulative percent. To calculate the percentile rank, half of 3.4% (1.7%) is added to 43.5% from Cumulative Percent to give you a percentile rank of 45.2, as shown in the table.

Table 1: Converting Tabulated Test Scores to NCE Values

Score	Frequency	Cumulative Frequency	Percent	Cumulative Percent	Percentile Rank	Z-Score	NCE
418	3,996	48,246	3.1	36.9	35.4	-0.375	42.10
420	4,265	52,511	3.3	40.2	38.5	-0.291	43.87
423	4,360	56,871	3.3	43.5	41.8	-0.206	45.66
425	4,404	61,275	3.4	46.9	45.2	-0.121	47.46
428	4,543	65,818	3.5	50.4	48.6	-0.035	49.27
430	4,619	70,437	3.5	53.9	52.1	0.053	51.12
432	4,645	75,082	3.6	57.4	55.7	0.143	53.00

NCEs are obtained from the percentile ranks using the normal distribution. The table of the standard normal distribution (found in many textbooks³) or computer software (for example, a spreadsheet) provides the associated Z-score from a standard normal distribution for any given percentile rank. NCEs are Z-scores that have been rescaled to have a “percentile-like” scale. As mentioned above, the NCE distribution is scaled so that NCEs exactly match the percentile ranks at 1, 50, and 99. To do this, each Z-score is multiplied by approximately 21.063 (the standard deviation on the NCE scale) and then 50 (the mean on the NCE scale) is added.

In previous years, the mCLASS assessment used book levels, which had a different process for converting to NCEs. The current Early Literacy assessment uses scale scores, so the process for converting to NCEs is similar to the EOG assessments.

With the test scores converted to NCEs, growth is calculated as the difference from one year and grade to the next in the same subject for a group of students. This process is explained in more technical detail in the next section.

2.2.3.3 How NCEs are Calculated for Non-Numeric Scales in Early Literacy Assessments

NCEs can also be created for assessments where the underlying scale is not inherently numeric in nature. One such assessment is the K-2 Text Reading and Comprehension assessment, which presents student achievement results in book levels and performance levels. Book levels range from Print Concepts (PC), Reading Behaviors (RB), B, C, and so on up to U. PC is the lowest possible book level and U is the highest possible book level on the distribution of possible book levels. Furthermore, each book level has three performance levels corresponding to the student’s reading and comprehension level of the text: Frustrational, Instructional, and Independent. Even though book levels and performance levels are non-numeric, the combination of the two provides the measured reading and comprehension ability of the test taker.

The frequencies of all observed book levels and performance levels of a population of test takers can be aggregated in an overall scoring distribution where each book and performance level are translated to corresponding percentiles and NCEs just as the case with other assessments that report numeric scale scores. NCEs for the K-2 Assessment in NC are calculated by grade and benchmark period: Beginning-of-Year (BOY), Middle-of-Year (MOY), and End-of-Year (EOY).

Growth for the Early Literacy assessments is the difference in NCEs from a starting benchmark period (MOY for Kindergartners, BOY for 1st and 2nd grade) to the EOY benchmark period. The average NCE for a district, school, or classroom can be compared to the overall amount of growth exhibited in the state, which represents a “normal” year’s growth, otherwise known as the growth standard.

³ See, for example, the inside front cover of William Mendenhall, Richard L. Scheaffer, and Dennis D. Wackerly, *Mathematical Statistics with Applications* (Boston: Duxbury Press, 1986).

2.2.4 Technical Description of the Gain Model

2.2.4.1 Definition of the Linear Mixed Model

As a linear mixed model, the gain model for district, school, and teacher value-added reporting is represented by the following equation in matrix notation:

$$y = X\beta + Zv + \epsilon \quad (1)$$

y (in the growth context) is the $m \times 1$ observation vector containing test scores (usually NCEs) for all students in all academic subjects tested over all grades and years.

X is a known $m \times p$ matrix that allows the inclusion of any fixed effects.

β is an unknown $p \times 1$ vector of fixed effects to be estimated from the data.

Z is a known $m \times q$ matrix that allows the inclusion of random effects.

v is a non-observable $q \times 1$ vector of random effects whose realized values are to be estimated from the data.

ϵ is a non-observable $m \times 1$ random vector variable representing unaccountable random variation.

Both v and ϵ have means of zero, that is, $E(v) = 0$ and $E(\epsilon) = 0$. Their joint variance is given by:

$$\text{Var} \begin{bmatrix} v \\ \epsilon \end{bmatrix} = \begin{bmatrix} G & 0 \\ 0 & R \end{bmatrix} \quad (2)$$

where R is the $m \times m$ matrix that reflects the amount of variation in and the correlation among the student scores residual to the specific model being fitted to the data, and G is the $q \times q$ variance-covariance matrix that reflects the amount of variation in and the correlation among the random effects. If (v, ϵ) are normally distributed, the joint density of (y, v) is maximized when β has value b and v has value u given by the solution to the following equations, known as Henderson's mixed model equations:⁴

$$\begin{bmatrix} X^T R^{-1} X & X^T R^{-1} Z \\ Z^T R^{-1} X & Z^T R^{-1} Z + G^{-1} \end{bmatrix} \begin{bmatrix} b \\ u \end{bmatrix} = \begin{bmatrix} X^T R^{-1} y \\ Z^T R^{-1} y \end{bmatrix} \quad (3)$$

Let a generalized inverse of the above coefficient matrix be denoted by

$$\begin{bmatrix} X^T R^{-1} X & X^T R^{-1} Z \\ Z^T R^{-1} X & Z^T R^{-1} Z + G^{-1} \end{bmatrix}^- = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = C \quad (4)$$

If G and R are known, then some of the properties of a solution for these equations are:

1. Equation (5) below provides the best linear unbiased estimator (BLUE) of the estimable linear function, $K^T \beta$, of the fixed effects. The second equation (6) below represents the variance of

⁴ McLean, Robert A., William L. Sanders, and Walter W. Stroup (1991). "A Unified Approach to Mixed Linear Models." *The American Statistician*, Vol. 45, No. 1, pp. 54-64.

that linear function. The standard error of the estimable linear function can be found by taking the square root of this quantity.

$$E(K^T \beta) = K^T b \quad (5)$$

$$\text{Var}(K^T b) = (K^T)C_{11}K \quad (6)$$

2. Equation (7) below provides the best linear unbiased predictor (BLUP) of v .

$$E(v|u) = u \quad (7)$$

$$\text{Var}(u - v) = C_{22} \quad (8)$$

where u is unique regardless of the rank of the coefficient matrix.

3. The BLUP of a linear combination of random and fixed effects can be given by equation (9) below provided that $K^T \beta$ is estimable. The variance of this linear combination is given by equation (10).

$$E(K^T \beta + M^T v | u) = K^T b + M^T u \quad (9)$$

$$\text{Var}(K^T(b - \beta) + M^T(u - v)) = (K^T M^T)C(K^T M^T)^T \quad (10)$$

4. With G and R known, the solution for the fixed effects is equivalent to generalized least squares, and if v and ϵ are multivariate normal, then the solutions for β and v are maximum likelihood.
5. If G and R are not known, then as the estimated G and R approach the true G and R , the solution approaches the maximum likelihood solution.
6. If v and ϵ are not multivariate normal, then the solution to the mixed model equations still provides the maximum correlation between v and u .

2.2.4.2 District and School Models

The district and school gain models do not contain random effects; consequently, the Zv term drops out in the linear mixed model. The X matrix is an incidence matrix (a matrix containing only zeros and ones) with a column representing each interaction of school (in the school model), subject, grade, and year of data. The fixed-effects vector β contains the mean score for each school, subject, grade, and year with each element of β corresponding to a column of X . Since gain models are generally run with each school uniquely defined across districts, there is no need to include districts in the model.

Unlike the case of the usual linear model used for regression and analysis of variance, the elements of ϵ are not independent. Their interdependence is captured by the variance-covariance matrix, which is also known as the R matrix. Specifically, scores belonging to the same student are correlated. If the scores in y are ordered so that scores belonging to the same student are adjacent to one another, then the R matrix is block diagonal with a block, R_i , for each student. Each student's R_i is a subset of the "generic" covariance matrix R_0 that contains a row and column for each subject and grade. Covariances among subjects and grades are assumed to be the same for all years (technically, all cohorts), but otherwise the R_0 matrix is unstructured. Each student's R_i contains only those rows and columns from R_0 that match

the subjects and grades for which the student has test scores. In this way, the gain model is able to use all available scores from each student.

Algebraically, the district gain model is represented as:

$$y_{ijkl} = \mu_{jkl} + \epsilon_{ijkl} \quad (11)$$

where y_{ijkl} represents the test score for the i^{th} student in the j^{th} subject in the k^{th} grade during the l^{th} year in the d^{th} district. μ_{jkl} is the estimated mean score for this particular district, subject, grade, and year. ϵ_{ijkl} is the random deviation of the i^{th} student's score from the district mean.

The school gain model is represented as:

$$y_{ijks} = \mu_{jks} + \epsilon_{ijks} \quad (12)$$

This is the same as the district analysis with the addition of the subscript s representing s^{th} school.

The gain model uses multiple years of student testing data to estimate the covariances that can be found in the matrix R_0 . This estimation of covariances is done within each level of analyses and can result in slightly different values within each analysis.

Solving the mixed model equations for the district or school gain model produces a vector b that contains the estimated mean score for each school (in the school model), subject, grade, and year. To obtain a value-added measure of average student growth, a series of computations can be done using the students from a school in a particular year and their prior and current testing data. The model produces means in each subject, grade, and year that can be used to calculate differences in order to obtain gains. Because students might change schools from one year to the next (in particular when transitioning from elementary to middle school, for example), the estimated mean score for the prior year/grade uses students who existed in the current year of that school. Therefore, mobility is taken into account within the model. Growth of students is computed using all students in each school including those that might have moved buildings from one year to the next.

The computation for obtaining a growth measure can be thought of as a linear combination of fixed effects from the model. The best linear unbiased estimate for this linear combination is given by equation (5). The growth measures are reported along with standard errors, and these can be obtained by taking the square root of equation (6) as described above.

2.2.4.3 Teacher Model

The teacher estimates use a more conservative statistical process to lessen the likelihood of misclassifying teachers. Each teacher's growth measure is assumed to be equal to either the state average or (for Early Literacy assessments) the average based on the population of test-takers in a specific year, subject, and grade until the weight of evidence pulls them either above or below that average. The model also accounts for the percentage of instructional responsibility the teacher has for each student during the course of each school year. Furthermore, the teacher model is "layered", which means that:

- Students' performance with both their current and previous teacher effects are incorporated.
- For each school year, the teacher estimates are based students' testing data collected over multiple previous years.

Each element of the statistical model for teacher value-added modeling provides an additional level of protection against misclassifying each teacher estimate.

To allow for the possibility of many teachers with relatively few students per teacher, the gain model enters teachers as random effects via the Z matrix in the linear mixed model. The X matrix contains a column for each subject, grade, and year, and the b vector contains an estimated state mean score for each subject, grade, and year. The Z matrix contains a column for each subject, grade, year, and teacher, and the u vector contains an estimated teacher effect for each subject, grade, year, and teacher. The R matrix is as described above for the district or school model. The G matrix contains teacher variance components with a separate unique variance component for each subject, grade, and year. To allow for the possibility that a teacher might be very effective in one subject and very ineffective in another, the G matrix is constrained to be a diagonal matrix. Consequently, the G matrix is a block diagonal matrix with a block for each subject/grade/year. Each block has the form $\sigma^2_{jkl}I$ where σ^2_{jkl} is the teacher variance component for the j^{th} subject in the k^{th} grade in the l^{th} year, and I is an identity matrix.

Algebraically, the teacher model is represented as:

$$y_{ijkl} = \mu_{jkl} + \left(\sum_{k^* \leq k} \sum_{t=1}^{T_{ijk^*l^*}} w_{ijk^*l^*t} \times \tau_{jk^*l^*t} \right) + \epsilon_{ijkl} \quad (13)$$

y_{ijkl} is the test score for the i^{th} student in the j^{th} subject in the k^{th} grade in the l^{th} year. $\tau_{jk^*l^*t}$ is the teacher effect of the t^{th} teacher in the j^{th} subject in grade k^* in year l^* . The complexity of the parenthesized term containing the teacher effects is due to two factors. First, in any given subject, grade, and year, a student might have more than one teacher. The inner (rightmost) summation is over all the teachers of the i^{th} student in a particular subject, grade, and year, denoted by $T_{ijk^*l^*}$. $\tau_{jk^*l^*t}$ is the effect of the t^{th} teacher. $w_{ijk^*l^*t}$ is the fraction of the i^{th} student's instructional time claimed by the t^{th} teacher. Second, as mentioned above, this model allows teacher effects to accumulate over time. The outer (leftmost) summation accumulates teacher effects not only for the current (subscripts k and l) but also over previous grades and years (subscripts k^* and l^*) in the same subject. Because of this accumulation of teacher effects, this type of model is often called the "layered" model.

In contrast to the model for many district and school estimates, the value-added estimates for teachers are not calculated by taking differences between estimated mean scores to obtain mean gains. Rather, this teacher model produces teacher "effects" (in the u vector of the linear mixed model). It also produces state-level mean scores (for each year, subject, and grade) in the fixed-effects vector b . Because of the way the X and Z matrices are encoded, in particular because of the "layering" in Z , teacher gains can be estimated by adding the teacher effect to the state mean gain. That is, the interpretation of a teacher effect in this teacher model is as a gain expressed as a deviation from the average gain for the state in a given year, subject, and grade.

Table 2 illustrates how the Z matrix is encoded for three students who have three different scenarios of teachers during grades 3, 4, and 5 in two subjects, Math (M) and Reading (R). Teachers are identified by the letters A–F, and students are identified by the letters X–Z.

Student X's teachers represent the conventional scenario. Student X is taught by a single teacher in both subjects each year (teachers A, C, and E in grades 3, 4, and 5, respectively). Notice that in Student X's Z matrix rows for grade 4 there are ones (representing the presence of a teacher effect) not only for fourth-grade teacher C but also for third-grade teacher A. This is how the "layering" is encoded. Similarly, in the grade 5 rows, there are ones for grade 5 teacher E, grade 4 teacher C, and grade 3 teacher A.

Student Y is taught by two different teachers in grade 3: teacher A for Math and teacher B for Reading. In grade 4, Student Y had teacher C for Reading. For some reason, in grade 4 no teacher claimed Student Y for Math even though Student Y had a grade 4 Math test score. This score can still be included in the analysis by entering zeros into the Student Y's Z matrix rows for grade 4 Math. In grade 5, however, Student Y had no test score in Reading. This row is completely omitted from the Z matrix. There will always be a Z matrix row corresponding to each test score in the y vector. Since Student Y has no entry in y for grade 5 Reading, there can be no corresponding row in Z .

Student Z's scenario illustrates team teaching. In grade 3 Reading, Student Z received an equal amount of instruction from teachers A and B. The entries in the Z matrix indicate each teacher's contribution, 0.5 for each teacher. In grade 5 Math, however, Student Z was taught by both teachers E and F, but they did not make an equal contribution. Teacher E claimed 80% responsibility, and teacher F claimed 20%.

Because teacher effects are treated as random effects in this approach, their estimates are obtained by shrinkage estimation, which is technically known as best linear unbiased prediction or as empirical Bayesian estimation. This means that *a priori* a teacher is considered "average" (with a teacher effect of zero) until there is sufficient student data to indicate otherwise. This method of estimation protects against false positives (teachers incorrectly evaluated as most effective or least effective), particularly in the case of teachers with few students.

Table 2: Encoding the Z Matrix

Student	Grade	Subjects	Teachers												
			Third Grade				Fourth Grade				Fifth Grade				
			A		B		C		D		E		F		
			M	R	M	R	M	R	M	R	M	R	M	R	
X	3	M	1	0	0	0	0	0	0	0	0	0	0	0	0
		R	0	1	0	0	0	0	0	0	0	0	0	0	0
	4	M	1	0	0	0	1	0	0	0	0	0	0	0	0
		R	0	1	0	0	0	1	0	0	0	0	0	0	0
	5	M	1	0	0	0	1	0	0	0	1	0	0	0	0
		R	0	1	0	0	0	1	0	0	0	1	0	0	0
Y	3	M	1	0	0	0	0	0	0	0	0	0	0	0	0
		R	0	0	0	1	0	0	0	0	0	0	0	0	0
	4	M	1	0	0	0	0	0	0	0	0	0	0	0	0
		R	0	0	0	1	0	1	0	0	0	0	0	0	0
	5	M	1	0	0	0	0	0	0	0	0	0	0	1	0
		R	0	0	0	0	0	0	0	0	0	0	0	0	0
Z	3	M	1	0	0	0	0	0	0	0	0	0	0	0	0
		R	0	0.5	0	0.5	0	0	0	0	0	0	0	0	0
	4	M	1	0	0	0	0	0	1	0	0	0	0	0	0
		R	0	0.5	0	0.5	0	0	0	1	0	0	0	0	0
	5	M	1	0	0	0	0	0	1	0	0.8	0	0.2	0	0
		R	0	0.5	0	0.5	0	0	0	1	0	0	0	0	1

From the computational perspective, the teacher gain can be defined as a linear combination of both fixed effects and random effects and is estimated by the model using equation (9). The variance and standard error can be found using equation (10).

2.2.4.4 Student Groups Model

The gain model provides district and school growth measures for their students included in a specific student group. In this analysis, expected growth is the same as in the overall students' analysis. In other words, expected growth is based on all students since the NCE mapping is based on all students, not just those in a specific student group. Furthermore, the estimated covariance parameters are used from the overall students' analysis when calculating the value-added measures.

Students are identified as members of a subgroup based on a flag in the student record. Growth measures are calculated for each subset of students for each district and school that meet the minimum requirements of student data.

2.2.4.5 Accommodations to the Gain Model in 2020-21 Reporting for Missing 2019-20 Data Due to the Pandemic

2.2.4.5.1 Overview

This section describes accommodations to the gain model that were made for *2020-21* reporting. However, this section does not apply to this year's reporting since the immediate prior year is available to measure student growth.

In spring 2020, the COVID-19 pandemic required schools to close early and cancel statewide summative assessments. As a result, scores are not available for North Carolina's EOG exams based on the 2019-20 school year, and it is not possible to measure growth on the EOGs from the 2018-19 to the 2019-20 school years or from the 2019-20 to the 2020-21 school years. For the gain model based on EOG exams, the 2020-21 reporting measures growth from the 2018-19 school year to the 2020-21 school year, except for EOG Reading for grade 3 and 4.

In grade 4, EVAAS measures growth from the beginning-of-year EOG Reading assessment in grade 3 from the 2018-19 school year (since there is no end-of-year score available) to the end-of-year EOG Reading assessment in grade 4 from the 2020-21 school year. In grade 3, EVAAS measures growth from the "beginning-of-year" EOG Reading assessment in grade 3 from the 2020-21 school year to the end-of-year EOG Reading assessment in grade 3 from the 2020-21 school year. It should be noted that the beginning-of-year assessments in grade 3 were not necessarily administered at the very beginning of year in each LEA. In fact, these assessments were administered through March 2021. To account for this administration window, SAS modified the NCEs attributed to each scale score on the beginning-of-year test based on when the student took the assessment as well as other available prior assessment data before grade 3. NCEs tended to be adjusted downward when the student took the assessment much later in the administration window. Any assessments that were administered after January 13th, 2021, were not used in this analysis due to investigations of the modeling approach and not being able to adequately estimate the entering achievement of these students.

From a technical perspective, the district and school gain model for EOGs is essentially the same as it has been in previous years except that growth is measured over two years rather than one year. Because

EVAAS measures the change in *relative* achievement based on the statewide population of test-takers, the growth measures are *relative* to the average growth observed in the state. In other words, a drop in achievement or proficiency rates due to lost instructional time does not correspond to a drop in growth. District, school, and teacher growth measures are still relative to the state average, and expected growth is based on students' maintaining their achievement among the population of test-takers.

That said, the interpretation of these growth measures changes slightly in two notable ways.

First, because the models provide two-year growth measures, the growth measure for grades where students transition from one school to another will then include growth from the feeder schools as well as the receiver school. For example, a middle school with grades 6–8 could receive a growth measure for sixth grade based on the students' growth in sixth grade as well as their growth from the feeder elementary schools in fifth grade.

In other words, it is not possible to parse out the individual contribution of the middle school in sixth grade apart from those from the elementary schools in fifth grade because of the missing year of test scores. For the district-level growth measures and for the non-transition grades, the two-year growth measures are still solely representative of growth within the specific district and the non-transition grades for the school are still solely representative of growth within the specific school.

Second, at a particular school, the growth of certain groups of students are not represented in the two-year measures as they would be in two one-year growth measures. For example, it is not possible to measure the growth of EOG Math for grade 4 this year because there is no EOG Math for grade 3 data from 2020. Similarly, it is not possible to report grade 8 growth from 2021 because there is no exiting achievement for these students in their last year at the school.

Despite these differences, the conceptual explanation of the 2020-21 growth measures is the same as it has always been: these growth measures compare students' exiting achievement with their entering achievement over two points in time.

Because Early Literacy assessments are administered at several points throughout the school year, the gain model for Early Literacy measures growth in the same way it has in previous years.

2.2.4.5.2 Additional Accommodations for the Teacher Model

The teacher-specific growth measure typically uses two types of information:

1. Student-level information, meaning the current and prior achievement data for the students connected to that teacher
2. Teacher-level information, meaning the present and past student-teacher linkages for the students connected to that teacher

The teacher gain model uses this information to estimate the current year teacher's contribution to the students' growth in the current year.

A teacher effect can be interpreted as the teacher's contribution to a student's gain. To estimate the *gain* for a teacher, the teacher *effect* is added to the statewide average gain. The statewide average gain is calculated in the same way as the gain in the district and school models. That is, in the absence of 2019-20 test scores, the 2020-21 statewide gain is a two-year gain. The teacher model differs from the

district and school models in using shrinkage estimation (“random” teacher effects) to provide estimates that are more stable with the smaller number of students typically associated with a teacher as compared with a school or district.

Even though there are no 2019-20 end-of-grade test scores, information about the 2019-20 teachers associated with each student is available, meaning the student-teacher linkage and percentage of instructional responsibility that the teacher has for each student in a given subject/grade in 2019-20. This information makes it possible to link the 2020-21 teachers *as well as the 2019-20 teachers* to the students with 2020-21 test scores. This information might improve estimates of the 2020-21 teacher measures by accounting for systematic differences observed in the student data when connected to specific teachers in the 2019-20 school year.

2.3 Predictive Model

2.3.1 Overview

Tests that are not given in consecutive grades require a different modeling approach from the gain model. The predictive model is used for such assessments in North Carolina. **The predictive model is a regression-based model where growth is a function of the difference between students’ expected scores with their actual scores.** Expected growth is met when students with a district, school, or teacher made the same amount of growth as students with the average district, school, or teacher.

Like the gain model, there are three separate analyses for EVAAS reporting based on the predictive model: one each for districts, schools, and teachers. The district and school models are essentially the same, and the teacher model includes accommodations for team teaching and other shared instruction.

Regression models are used in virtually every field of study, and their intent is to identify relationships between two or more variables. When it comes to measuring growth, regression models identify the relationship between prior test performance and actual test performance for a given course. In more technical terms, the predictive model is known as the univariate response model (URM), a linear mixed model and, more specifically, an analysis of covariance (ANCOVA) model.

The key advantages of the predictive model can be summarized as follows:

- It minimizes the influence of measurement error and increases the precision of predictions by using multiple prior test scores as predictors for each student.
- It does not require students to have all predictors or the same set of predictors as long as a student has at least three prior test scores as predictors of the response variable in any subject and grade.
- It allows educators to benefit from all tests, even when tests are on differing scales.
- It accommodates teaching scenarios where more than one teacher has responsibility for a student’s learning in a specific subject, grade, and year.

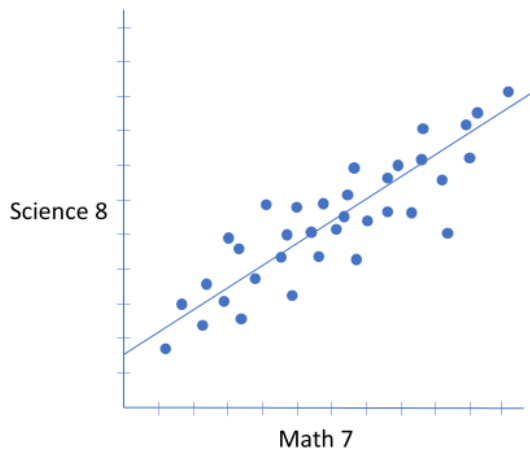
2.3.2 Conceptual Explanation

As mentioned above, the predictive model is ideal for assessments given in non-consecutive grades, such as EOG Science tests in grades 5 and 8, or the high school end-of-course tests, such as Math 1. It is also used for EOG Mathematics and Reading in grade 8 due to the prevalence of accelerated testing for

grade 8 Math students, which means that the grade 8 cohort has a non-trivial number of students who no longer test in consecutive grades in grade 8.

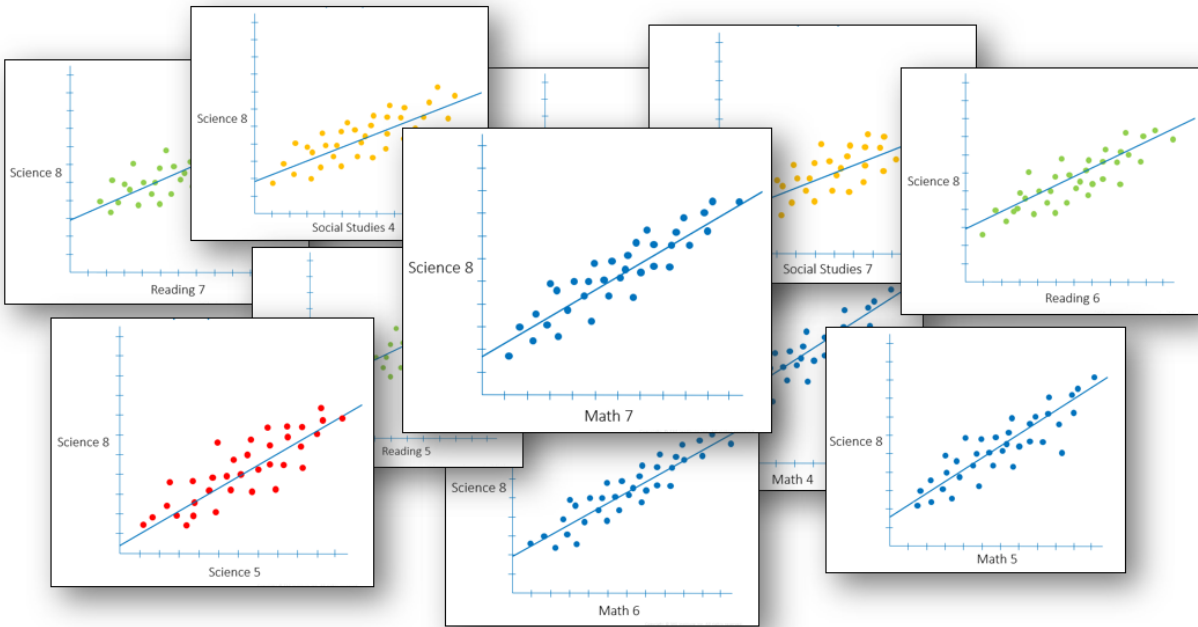
Consider all students who tested in EOG Science in grade 8 in a given year. The gain model is not possible since there isn't a Science test in the immediate prior grade. However, these students might have a number of prior test scores in EOG Math and Reading in grades 3–7 and EOG Science in grade 5. These prior test scores have a relationship with EOG Science, meaning that how students performed on these tests can predict how the students perform on EOG Science in grade 8. The growth model does not assume what the predictive relationship will be; instead, the actual relationships observed by the data define the relationships. This is shown in Figure 3 below where each dot represents a student's prior score on EOG Math 7 plotted with their score on EOG Science 8. The best-fit line indicates how students with a certain prior score on EOG Math 7 tend to score, on average, on EOG Science 8. This illustration is based on one prior test; the predictive model uses many prior test scores from different subjects and grades.

Figure 3: Test Scores from One Assessment Have a Predictive Relationship to Test Scores from Another Assessment



Some subjects and grades will have a greater relationship to EOG Science in grade 8 than others; however, the other subjects and grades still have a predictive relationship. For example, prior Math scores might have a stronger predictive relationship to EOG Science in grade 8 than prior Reading scores, but how a student reads and performs on the EOG Reading test typically provides an idea of how we might expect a student to perform on average on EOG Science. This is shown in Figure 4 below where there are a number of different tests that have a predictive relationship with EOG Science in grade 8. All of these relationships are considered together in the predictive model with some tests weighted more heavily than others.

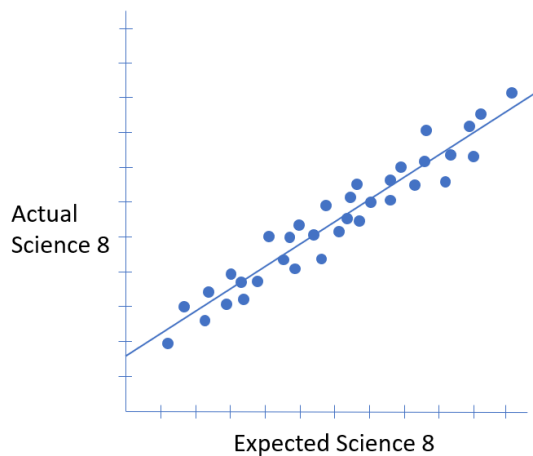
Figure 4: Relationships Observed in the Statewide Data Inform the Predictive Model



Note that the prior test scores do not need to be on the same scale as the assessment being measured for student growth. Just as height (reported in inches) and weight (reported in pounds) can predict a child's age (reported in years), the growth model can use test scores from different scales to find the predictive relationship.

Each student receives an expected score based on their own prior testing history. In practical terms, the expected score represents the student's entering achievement because it is based on all prior testing information to date. Figure 5 below shows the relationship between expected and actual scores for a group of students.

Figure 5: Relationship Expected Score and Actual Score for Selected Subject and Grade



The expected scores can be aggregated to a specific district, school, or teacher and then compared to the students' actual scores. In other words, the growth measure is a function of the difference between the exiting achievement (or average actual score) and the entering achievement (or average expected score) for a group of students. Unlike the gain model, the actual score and expected score are reported in the scaling units of the test rather than NCEs.

2.3.3 Technical Description of the District, School, and Teacher Models

The predictive model has similar approaches for districts and schools and a slightly different approach for teachers that accounts for shared instructional responsibility. The approach is described briefly below with more details following.

- The score to be predicted serves as the response variable (y , the dependent variable).
- The covariates (x terms, predictor variables, explanatory variables, independent variables) are scores on tests the student has taken in previous years from the response variable.
- There is a categorical variable (class variable, grouping variable) to identify the district, school, or teachers from whom the student received instruction in the subject, grade, and year of the response variable (y).

Algebraically, the model can be represented as follows for the i^{th} student, assuming in the teacher model that there is no team teaching.

$$y_i = \mu_y + \alpha_j + \beta_1(x_{i1} - \mu_1) + \beta_2(x_{i2} - \mu_2) + \dots + \epsilon_i \quad (14)$$

In the case of team teaching, the single α_j is replaced by multiple α s, each multiplied by an appropriate weight, similar to the way this is handled in the teacher gain model in equation (13). The μ terms are means for the response and the predictor variables. α_j is the teacher effect for the j^{th} district, school, or teacher—the one who claimed responsibility for the i^{th} student. The β terms are regression coefficients. Predictions to the response variable are made by using this equation with estimates for the unknown parameters (μ terms, β terms, and sometimes α_j). The parameter estimates (denoted with "hats," e.g., $\hat{\mu}$, $\hat{\beta}$) are obtained using all students that have an observed value for the specific response and have three predictor scores. The resulting prediction equation for the i^{th} student is as follows:

$$\hat{y}_i = \hat{\mu}_y + \hat{\beta}_1(x_{i1} - \hat{\mu}_1) + \hat{\beta}_2(x_{i2} - \hat{\mu}_2) + \dots \quad (15)$$

Two difficulties must be addressed in order to implement the prediction model. First, not all students will have the same set of predictor variables due to missing test scores. Second, because the predictive model is an ANCOVA model, the estimated parameters are pooled within group (district, school, or teacher). The strategy for dealing with missing predictors is to estimate the joint covariance matrix (call it C) of the response and the predictors. Let C be partitioned into response (y) and predictor (x) partitions, that is,

$$C = \begin{bmatrix} c_{yy} & c_{yx} \\ c_{xy} & c_{xx} \end{bmatrix} \quad (16)$$

Note that C in equation (16) is not the same as C in equation (4). This matrix is estimated using the EM (expectation maximization) algorithm for estimating covariance matrices in the presence of missing data available in SAS/STAT® (although no imputation is actually used). It should also be noted that, due to this being an ANCOVA model, C is a pooled-within group (district, school, or teacher) covariance matrix. This

is accomplished by providing scores to the EM algorithm that are centered around group means (i.e., the group means are subtracted from the scores) rather than around grand means. Obtaining C is an iterative process since group means are estimated within the EM algorithm to accommodate missing data. Once new group means are obtained, another set of scores is fed into the EM algorithm again until C converges. This overall iterative EM algorithm is what accommodates the two difficulties mentioned above. Only students who had a test score for the response variable in the most recent year and who had at least three predictor variables are included in the estimation. Given such a matrix, the vector of estimated regression coefficients for the projection equation (15) can be obtained as:

$$\hat{\beta} = C_{xx}^{-1} c_{xy} \quad (17)$$

This allows one to use whichever predictors a student has to get that student's expected y -value (\hat{y}_i). Specifically, the C_{xx} matrix used to obtain the regression coefficients *for a particular student* is that subset of the overall C matrix that corresponds to the set of predictors for which this student has scores.

The prediction equation also requires estimated mean scores for the response and for each predictor (the $\hat{\mu}$ terms in the prediction equation). These are not simply the grand mean scores. It can be shown that in an ANCOVA if one imposes the restriction that the estimated "group" effects should sum to zero (that is, the effect for the "average" district, school or teacher is zero), then the appropriate means are the means of the group means. The group-level means are obtained from the EM algorithm mentioned above, which accounts for missing data. The overall means ($\hat{\mu}$ terms) are then obtained as the simple average of the group-level means.

Once the parameter estimates for the prediction equation have been obtained, predictions can be made for any student with any set of predictor values as long as that student has a minimum of three prior test scores. This is to avoid bias due to measurement error in the predictors.

$$\hat{y}_i = \hat{\mu}_y + \hat{\beta}_1(x_{i1} - \hat{\mu}_1) + \hat{\beta}_2(x_{i2} - \hat{\mu}_2) + \dots \quad (18)$$

The \hat{y}_i term is nothing more than a composite of all the student's past scores. It is a one-number summary of the student's level of achievement prior to the current year, and this term is called the expected score or entering achievement in the web reporting. The different prior test scores making up this composite are given different weights (by the regression coefficients, the $\hat{\beta}$ terms) in order to maximize its correlation with the response variable. Thus, a different composite would be used when the response variable is Math than when it is Reading, for example. Note that the $\hat{\alpha}_j$ term is not included in the equation. Again, this is because \hat{y}_i represents prior achievement before the effect of the current district, school, or teacher.

The second step in the predictive model is to estimate the group effects (α_j) using the following ANCOVA model.

$$y_i = \gamma_0 + \gamma_1 \hat{y}_i + \alpha_j + \epsilon_i \quad (19)$$

In the predictive model, the effects (α_j) are considered random effects. Consequently, the $\hat{\alpha}_j$ terms are obtained by shrinkage estimation (empirical Bayes).⁵ The regression coefficients for the ANCOVA model are given by the γ terms.

2.3.3.1 Accommodations to the Predictive Model in 2020-21 Reporting for Missing 2019-20 Data due to the Pandemic

This section describes how the predictive model was calculated for 2020-21 reporting. However, this section does not apply to this year's reporting since the immediate prior year is available to measure student growth.

In spring 2020, the COVID-19 pandemic required schools to close early and cancel statewide summative assessments. As a result, scores are not available for North Carolina's EOG and EOC exams based on the 2019-20 school year, and it is not possible to measure growth from the 2018-19 to the 2019-20 school years. For the predictive model, the 2020-21 reporting measures growth using students' predictors through the 2019-20 school year where available and then compares to their performance on the 2020-21 assessment.

As a reminder, the predictive model is used to measure growth for assessments given in non-consecutive grades, such as EOG Science in grade 8 as well as the EOC assessments in Biology, English II, and Mathematics 1 and 3. Because these assessments are not administered every year, *it has always been possible* that students do not have any test scores in the *immediate* prior year. The model can provide a robust estimate of students' entering achievement for the course by using all other available test scores from other subjects, grades, and years.

In other words, the predictive model does not require any technical adaptations to account for the missing year of data and the interpretation of the results is similar to a typical year of reporting.

Unlike the gain-based model, the teacher predictive model only uses the current year teacher within the model. Given this, the teacher model was run in similar ways as in prior years.

2.4 Projection Model

2.4.1 Overview

The longitudinal data sets used to calculate growth measures for groups of students can also provide individual student projections to future assessments. A projection is reported as a probability of obtaining a specific score or above on an assessment, such as a 70% probability of scoring Level 3 or above on the next summative assessment. The probabilities are based on the students' own prior testing history as well as how the cohort of students who just took the assessment performed. Projections are available for state assessments as well as to college readiness assessments.

⁵ For more information about shrinkage estimation, see, for example, Ramon C. Littell, George A. Milliken, Walter W. Stroup, Russell D. Wolfinger, and Oliver Schabenberger, *SAS for Mixed Models, Second Edition* (Cary, NC: SAS Institute Inc., 2006). Another example is Charles E. McCulloch, Shayle R. Searle, and John M. Neuhaus, *Generalized, Linear, and Mixed Models, Second Edition* (Hoboken, NJ: John Wiley & Sons, 2008).

Projections are useful as a planning resource for educators, and they can inform decisions around enrollment, enrichment, remediation, counseling, and intervention to increase students' likelihood of future success.

2.4.2 Technical Description

The statistical model that is used as the basis for the projections is, in traditional terminology, an analysis of covariance (ANCOVA) model. This model is the same statistical model used in the predictive model applied at the school level described in [Section 2.3.3](#). In the projection model, the score to be projected serves as the response variable (y), the covariates (x terms) are scores on tests the student has already taken, and the categorical variable is the school at which the student received instruction in the subject, grade, and year of the response variable (y). Algebraically, the model can be represented as follows for the i^{th} student.

$$y_i = \mu_y + \alpha_j + \beta_1(x_{i1} - \mu_1) + \beta_2(x_{i2} - \mu_2) + \dots + \epsilon_i \quad (20)$$

The μ terms are means for the response and the predictor variables. α_j is the school effect for the j^{th} school, the school attended by the i^{th} student. The β terms are regression coefficients. Projections to the future are made by using this equation with estimates for the unknown parameters (μ terms, β terms, sometimes α_j). The parameter estimates (denoted with "hats," e.g., $\hat{\mu}$, $\hat{\beta}$) are obtained using the most current data for which response values are available. The resulting projection equation for the i^{th} student is

$$\hat{y}_i = \hat{\mu}_y \pm \hat{\alpha}_j + \hat{\beta}_1(x_{i1} - \hat{\mu}_1) + \hat{\beta}_2(x_{i2} - \hat{\mu}_2) + \dots \quad (21)$$

The reason for the " \pm " before the $\hat{\alpha}_j$ term is that since the projection is to a future time, the school that the student will attend is unknown, so this term is usually omitted from the projections. This is equivalent to setting $\hat{\alpha}_j$ to zero, that is, to assuming that the student encounters the "average schooling experience" in the future.

Two difficulties must be addressed to implement the projections. First, not all students will have the same set of predictor variables due to missing test scores. Second, because this is an ANCOVA model with a school effect i , the regression coefficients must be "pooled-within-school" regression coefficients. The strategy for dealing with these difficulties is the same as described in Section 2.3.3 using equations (16), (17), and (18) and will not be repeated here.

The parameter estimates are based on the cohort of students who most recently took the assessment, which is the 2022-23 school year for this year's reporting. Once the parameter estimates for the projection equation have been obtained, projections can be made for any student with any set of predictor values. However, to protect against bias due to measurement error in the predictors, projections are made only for students who have at least three available predictor scores. In addition to the projected score itself, the standard error of the projection is calculated ($SE(\hat{y}_i)$). Given a projected score and its standard error, it is possible to calculate the probability that a student will reach some specified benchmark of interest (b). Examples are the probability of scoring at least Proficient on a future end-of-grade test or the probability of scoring at least an established college readiness benchmark. The probability is calculated as the area above the benchmark cutoff score using a normal distribution with its mean equal to the projected score and its standard deviation equal to the standard

error of the projected score as described below. Φ represents the standard normal cumulative distribution function.

$$Prob(\hat{y}_i \geq b) = \Phi\left(\frac{\hat{y}_i - b}{SE(\hat{y}_i)}\right) \quad (22)$$

2.5 Outputs from the Models

2.5.1 Gain Model

The gain model is used for courses where students test in consecutive grade-given tests. As such, **the gain model uses EOG in Math and Reading in grades 3–8 and Early Literacy assessments in grades K-2 to provide district, school, and teacher growth measures in the following content areas:**

- EOG Math in grades 4–7 (included grade 8 for 2021-22 and 2020-21 reporting)
- EOG Reading in grades 3–7 (included grade 8 for 2021-22 and 2020-21 reporting)
- Early Literacy in grades K-2 based on State Board of Education-approved assessments (mCLASS for 2022-23 and 2021-22 reporting; i-Ready, Istation Indicators of Progress, mCLASS, NWEA Measures of Academic Progress (MAP) and Renaissance Learning Standardized Test for the Assessment of Reading (STAR) for 2020-21 reporting)

In addition to the mean scores and mean gain for an individual subject, grade, and year, the gain model can also provide the following:

- Cumulative gains across grades (for each subject and year)
- Multi-year up to 3-average gains (for each subject and grade) (not available for 2022-23 reporting)
- Composite gains across subjects

In general, these are all different forms of linear combinations of the fixed effects (and random effects for the teacher model), and their estimates and standard errors are computed in the same manner described above in equations (5) and (6) for district and school models and in equations (9) and (10) for the teacher model.

Collectively, the different models provide metrics for a variety of purposes within the State of North Carolina. They are summarized in the table below, and note that Early Literacy reports only for the overall students measures while EOGs report measures for overall students as well as the student groups where sufficient data exists:

- District growth measures
 - Overall students
 - Academically or Intellectually Gifted (AIG)
 - American Indian/Alaskan Native
 - Asian/Pacific Islander
 - Black (not Hispanic)
 - Economically Disadvantaged Students (EDS)
 - English Learners (EL)
 - Hispanic

- Homeless
- Students with Disabilities (SWD)
- Two or More Races
- White (not Hispanic)
- School growth measures
 - Overall students
 - Academically or Intellectually Gifted (AIG)
 - American Indian/Alaskan Native
 - Asian/Pacific Islander
 - Black (not Hispanic)
 - Economically Disadvantaged Students (EDS)
 - English Learners (EL)
 - Hispanic
 - Homeless
 - Students with Disabilities (SWD)
 - Two or More Races
 - White (not Hispanic)
- Teacher growth measures based on linked students

More details about district, school, and teacher composites across subjects, grades, and years are available in [Section 5](#).

2.5.2 Predictive Model

The predictive model is used for courses where students test in non-consecutive grade-given tests. As such, **the predictive model provides growth measures for districts, schools, and teachers in the following content areas:**

- EOG Mathematics and Reading in grade 8
- EOG Science in grade 5 and 8
- EOC Biology
- EOC English II
- EOC Math 1 and 3
 - Note: There is also an EOC Math 3 School Accountability Growth (SAG) measure created for students who took the Math 1 assessment in middle school. This measure is only available at the school level.

The predictive model also provides district- and school-level growth measures only in the following content areas:

- ACT English, Math, Reading, Science, and Composite
- CTE for a variety of content areas depending on the availability of CTE courses and sufficient participation among students, schools, and districts
- PSAT NMSQT Evidence-Based Reading and Writing, Mathematics, and Total Score
- SAT Evidence-Based Reading and Writing, Mathematics, and Total Score

In addition to the mean scores and growth measures for an individual subject, grade, and year, the predictive model can also provide multi-year average growth measures (up to three years) for each

subject and grade or course. This multi-year average is not available for 2022-23 reporting due to the missing year of data.

Collectively, the different models provide metrics for a variety of purposes within the State of North Carolina. They are summarized in the table below. Note that EOG Science, Math 1, Math 3, and English II report measures for overall students as well as the student groups where sufficient data exists :

- District growth measures
 - Overall students
 - Academically or Intellectually Gifted (AIG)
 - American Indian/Alaskan Native
 - Asian/Pacific Islander
 - Black (not Hispanic)
 - Economically Disadvantaged Students (EDS)
 - English Learners (EL)
 - Hispanic
 - Homeless
 - Students with Disabilities (SWD)
 - Two or More Races
 - White (not Hispanic)
- School growth measures
 - Overall students
 - Academically or Intellectually Gifted (AIG)
 - American Indian/Alaskan Native
 - Asian/Pacific Islander
 - Black (not Hispanic)
 - Economically Disadvantaged Students (EDS)
 - English Learners (EL)
 - Hispanic
 - Homeless
 - Students with Disabilities (SWD)
 - Two or More Races
 - White (not Hispanic)
- Teacher growth measures based on linked students

More details about district, school, and teacher composites across subjects, grades, and years are available in [Section 5](#).

2.5.3 Projection Model

Projections are provided to future state assessments as well as college readiness assessments. More specifically, EOG projections are typically provided to a student's next two tested grade-level EOG assessments based on that student's most recent tested grade, such as projections to grades 6 and 7 for students who most recently tested in grade 5. For the 2022-23 reporting, EOG projections are provided to a student's next tested grade-level EOG assessment, such as a projection to grade 6 for students who most recently tested in grade 5. EOC projections are provided for students as soon as they have at least three test scores in common with the students in the most recent tested cohort. Projections are made to the performance levels 3– 5, depending on the assessment, and the individual cut scores depend on

each subject and grade. Due to the lack of prior test scores, projections for grades K-3 are not available for the 2022-23 reporting.

CTE projections are available in many subject areas, and the exact offering varies each year depending on the availability of CTE courses and sufficient numbers of students taking the assessment. The CTE projections are made to the Proficient level for each course.

ACT, PSAT, or SAT projections are provided to students who last tested in grades 4–11. ACT, PSAT, and SAT projections will be provided for the following subject areas:

- ACT Composite
- ACT English
- ACT Mathematics
- ACT Reading
- ACT Science
- PSAT NMSQT Evidence-Based Reading and Writing
- PSAT NMSQT Mathematics
- SAT Total Score
- SAT Evidence-Based Reading and Writing
- SAT Mathematics

Advanced Placement (AP) projections are provided in the following subject areas:

- AP Biology
- AP Calculus AB
- AP Calculus BC
- AP Computer Science Principles
- AP English Language
- AP English Literature
- AP Environmental Science
- AP Human Geography
- AP Psychology
- AP Statistics
- AP United States Government and Politics
- AP United States History
- AP World History

3 Expected Growth

3.1 Overview

Conceptually, growth is simply the difference between students' entering and exiting achievement. As noted in Section 2, zero represents "expected growth." Positive growth measures are evidence that students made *more* than the expected growth, and negative growth measures are evidence that students made *less* than the expected growth.

A more detailed explanation of expected growth and how it is calculated are useful for the interpretation and application of growth measures.

3.2 Technical Description

Both the gain and predictive models define expected growth based on the empirical student testing data; in other words, the model does not assume a particular amount of growth or assign expected growth in advance of the assessment being taken by students. Both models define expected growth within a year. This means that expected growth is always relative to how students' achievement has changed in the most recent year of testing rather than a fixed year in the past.

More specifically, **in the gain model, expected growth means that students maintained the same relative position with respect to the statewide student achievement that year. In the predictive model, expected growth means that students with a district, school, or teacher made the same amount of growth as students with the average district, school, or teacher in the state for that same year, subject, and grade.**

For both models, the growth measures tend to be centered on expected growth every year with approximately half of the district/school/teacher estimates above zero and approximately half of the district/school/teacher estimates below zero.

A change in assessments or scales from one year to the next does not present challenges to calculating expected growth. Through the use of NCEs, the gain model converts any scale to a relative position, and the predictive model already uses prior test scores from different scales to calculate the expected score. When assessments change over time, expected growth is still based on the relative change in achievement from one point in time to another.

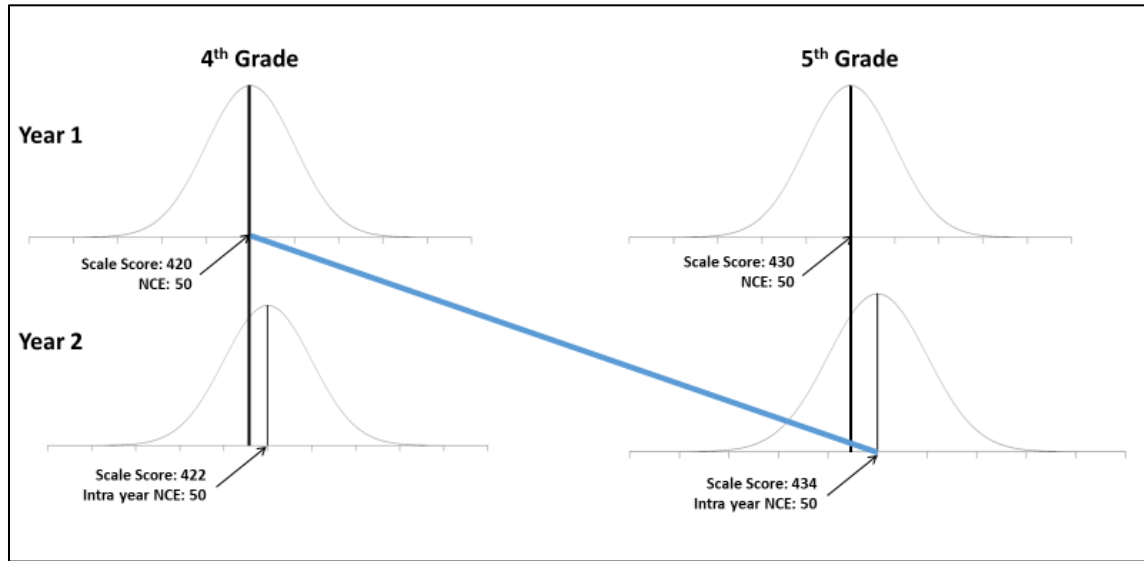
3.3 Illustrated Example

Figure 6 below provides a simplified example of how growth is calculated in the gain model when the state achievement increases. The figure has four graphs, each of which plot the NCE distribution of scale scores for a given year and grade. In this example, the figure shows how the gain is calculated for a group of grade 4 students in Year 1 as they become grade 5 students in Year 2. In Year 1, our grade 4 students score, on average, 420 scale score points on the test, which corresponds to the 50th NCE (similar to the 50th percentile). In Year 2, the students score, on average, 434 scale score points on the test, which corresponds to a 50th NCE *based on the grade 5 distribution of scores in Year 2*. The grade 5 distribution of scale scores in Year 2 was higher than the grade 5 distribution of scale scores in Year 1, which is why the lower right graph is shifted slightly to the right. The blue line shows what is required for students to make expected growth, which would be to maintain their position at the 50th NCE for grade

4 in Year 1 as they become grade 5 students in Year 2. The growth measure for these students is Year 2 NCE – Year 1 NCE, which would be 50 – 50 = 0. Similarly, if a group of students started at the 35th NCE, the expectation is that they would maintain that 35th NCE.

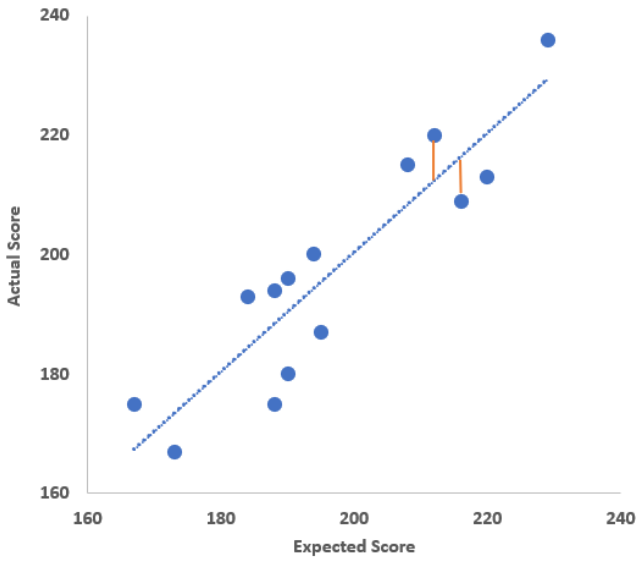
Note that the actual gain calculations are much more robust than what is presented here; as described in the previous section, the models can address students with missing data, team teaching, and all available testing history.

Figure 6: Intra-Year Approach Example for the Gain Model



In contrast, in the predictive model, expected growth uses actual results from the most recent year of assessment data and considers the relationships from the most recent year with prior assessment results. Figure 7 below provides a simplified example of how growth is calculated in the predictive model. The graph plots each student’s actual score with their expected score. Each dot represents a student, and a best-fit line will minimize the difference between all students’ actual and expected scores. Collectively, the best-fit line indicates what expected growth is for each student – given the student’s expected score, expected growth is met if the student scores the corresponding point on the best-fit line. Conceptually, with the best-fit line minimizing the difference between all students’ actual and expected scores, the growth expectation is defined by the average experience. Note that the actual calculations differ slightly since this is an ANCOVA model where the students are expected to see the average growth as seen by the experience with the average group (district, school, or teacher).

Figure 7: Intra-Year Approach Example for the Predictive Model



4 Classifying Growth into Categories

4.1 Overview

It can be helpful to classify growth into different levels for interpretation and context, particularly when the levels have statistical meaning. North Carolina's growth model has three categories for districts, schools, and teachers. These categories are defined by a range of values related to the growth measure and its standard error, and they are known as growth indicators in the web application.

4.2 Use Standard Errors Derived from the Models

As described in the modeling approaches section, the growth model provides an estimate of growth for a district, school, or teacher in a particular subject, grade, and year as well as that estimate's standard error. The standard error is a measure of the quantity and quality of student data included in the estimate, such as the number of students and the occurrence of missing data for those students. It also takes into account shared instruction and team teaching. Standard error is a common statistical metric reported in many analyses and research studies because it yields important information for interpreting an estimate, in this case the growth measure relative to expected growth. Because measurement error is inherent in any growth or value-added model, *the standard error is a critical part of the reporting.* **Taken together, the growth measure and standard error provide educators and policymakers with critical information about the certainty that students in a district, school, or classroom are making decidedly more or less than the expected growth.** Taking the standard error into account is particularly important for reducing the risk of misclassification (for example, identifying a teacher as ineffective when they are truly effective) for high-stakes usage of value-added reporting.

The standard error also takes into account that even among teachers with the same number of students, teachers might have students with very different amounts of prior testing history. Due to this variation, the standard errors in a given subject, grade, and year could vary significantly among teachers, depending on the available data that is associated with their students, and it is another important protection for districts, schools, and teachers to incorporate standard errors to the value-added reporting.

4.3 Define Growth Indicators in Terms of Standard Errors

Common statistical usage of standard errors indicates the precision of an estimate and whether that estimate is statistically significantly different from an expected value. The growth reports use the standard error of each growth measure to determine the statistical evidence that the growth measure is different from expected growth. For EVAAS growth reporting, this is essentially when the growth measure is more than or less than two standard errors above or below expected growth or, in other words, when the growth index is more than +2 or less than -2. These definitions then map to growth indicators in the reports themselves, such that there is statistical meaning in these categories. The categories and definitions are illustrated in the following section.

4.4 Illustrated Examples of Categories

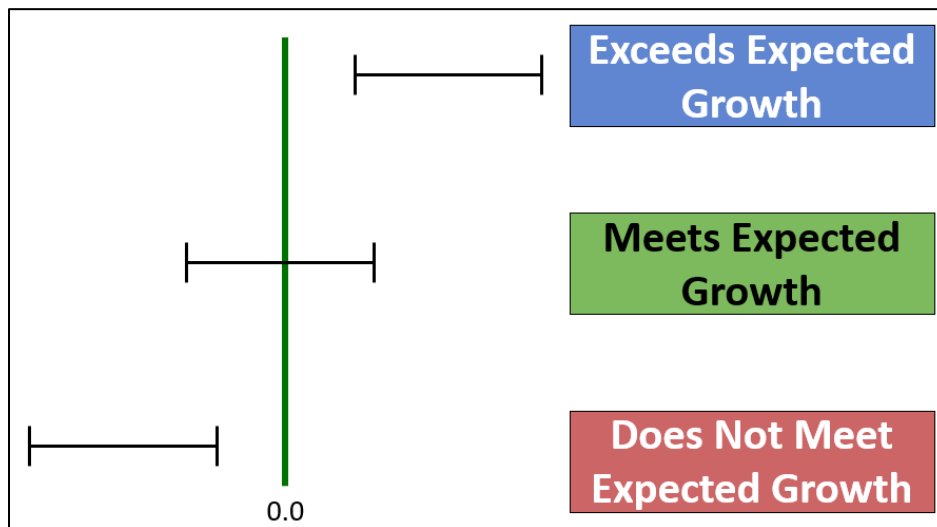
There are two ways to visualize how the growth measure and standard error relate to expected growth and how these can be used to create categories.

The first way is to frame the growth measure relative to its standard error and expected growth at the same time. For district and school reporting, the categories are defined as follows:

- **Exceeds Expected Growth** indicates that the growth measure is two standard errors or more above expected growth (0). This level of certainty is significant evidence of exceeding the standard for academic growth.
- **Meets Expected Growth** indicates that the growth measure is less than two standard errors above expected growth (0) and no more than two standard errors below it (0). This is evidence of meeting the standard for academic growth.
- **Does Not Meet Expected Growth** is an indication that the growth measure is more than two standard errors below expected growth (0). This level of certainty is significant evidence of not meeting the standard for academic growth.

Figure 8 below shows visual examples of each category. The green line represents the expected growth. The black line extends the range of values to the growth measure plus and minus two standard errors. If the black line is completely above expected growth, then there is significant evidence that students made more than expected growth, which represents the Exceeds Expected Growth category. Conversely, if the black line is completely below expected growth, then there is significant evidence that students made less than expected growth, which represents the Does Not Meet Expected Growth category. Meets Expected Growth indicates that there is evidence that students made growth as expected as the black line crosses the green line indicating expected growth.

Figure 8: Visualization of Growth Categories with Expected Growth, Growth Measures, and Standard Errors



This graphic is helpful in understanding how the growth measure relates to expected growth and whether the growth measure represents a statistically significant difference from expected growth.

The second way to illustrate the categories is to create a growth index, which is calculated as shown below:

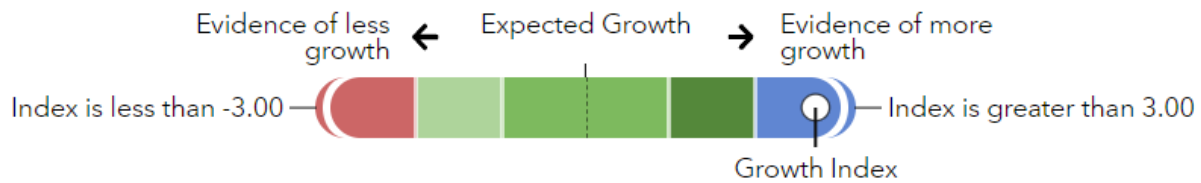
$$\text{Growth Index} = \frac{\text{Growth Measure} - \text{Expected Growth}}{\text{Standard Error of the Growth Measure}} \quad (23)$$

The growth index is similar in concept to a Z-score or t-value, and it communicates as a single metric the certainty or evidence that the growth measure is decidedly above or below expected growth. The growth index is useful when comparing value-added measures from different assessments or in different units, such as NCEs or scale scores. The categories can be established as ranges based on the growth index, such as the following:

- **Exceeds Expected Growth (Dark Blue)** indicates significant evidence that students made more growth than expected. The growth index is 2 or greater.
- **Meets Expected Growth (Green)** indicates evidence that students made growth as expected. The growth index is between -2 and 2.
- **Does Not Meet Expected Growth (Red)** indicates significant evidence that students made more growth than expected. The growth index is less than -2.

This is represented in the growth indicator bar in Figure 9, which is similar to what is provided in the District and School Value-Added reports in the EVAAS web application. The black dotted line represents expected growth. The color-coding within the bar indicates the range of values for the growth index within each category.

Figure 9: Sample Growth Indicator Bar



It is important to note that these two illustrations provide users with the same information; they are simply presenting the growth measure, its standard error, and expected growth in different ways.

4.5 Rounding and Truncating Rules

As described in the previous section, the effectiveness level is based on the value of the growth index. As additional clarification, the calculation of the growth index uses unrounded values for the value-added measures and standard errors. After the growth index has been created but before the categories are determined, the index values are rounded or truncated by taking the maximum value of the rounded or truncated index value out to two decimal places. This provides the highest category given any type of rounding or truncating situation. For example, if the score was a 1.995, then rounding would provide a higher category. If the score was a -2.005, then truncating would provide a higher category. In practical terms, this impacts only a very small number of measures.

Also, when value-added measures are combined to form composites, as described in the next section, the rounding or truncating occurs after the final index is calculated for that combined measure.

5 Composite Growth Measures

A composite combines growth measures from different subjects, grades, and/or courses. The key policy decisions for combining growth measures can be summarized as follows:

- The composite is based on the most recent year of growth measures.
- For schools, there is a School Accountability Growth (SAG) composite that only includes EOC and EOG subjects and grades. Beginning with the 2017-18 school year, Biology is not included in the School Accountability Growth composite.
- For schools and teachers, there is a second composite, which uses more subjects and grades associated with a school since it also includes Early Literacy and CTEs. This measure is called Educator Effectiveness Growth (EEG) for schools, and Student Growth Measure (SGM) for teachers.
- The SAG, EEG, and SGM composites weigh each subject/grade (for EOG and Early Literacy) and each subject (for EOC and CTE) according to the number of scores included in the growth measures.

The following sections show how an SGM composite is calculated for a sample teacher. Although we present a teacher example, the process for school composite calculations is the same for the EEG.

5.1 Teacher Composites

5.1.1 Overview

The key steps for determining a teacher's composite index are as follows:

1. For growth measures based on the gain model, calculate composite index across subjects and years.
2. For growth measures based on the predictive model, calculate composite index across subjects and years.
3. Using both the gain and predictive model composite indices, calculate the composite index.

If a teacher does not have value-added measures from both the gain and predictive models, then the composite index would be based on the model for which the teacher does have reporting. For the 2022-23 reporting, the composite is a single-year measure.

The following sections illustrate this process using value-added measures from a sample teacher, which are provided in Table 3.

Table 3: Sample Teacher Value-Added Information

Year	Subject	Grade	Value-Added Measure	Standard Error	Number of FYE Students
1	EOG Reading	7	-0.30	1.20	65
1	EOG Math	7	3.80	1.50	70
1	Math I	8	11.75	6.20	20

5.1.2 Technical Description of the Composite Index Based on Gain Model Measures

The composite index for the gain model growth measures is calculated by dividing the composite gain by its composite standard error. The calculations for each of these metrics are provided below.

5.1.2.1 Composite Gain Across Subjects

Growth measures from the gain model are in the same scale (NCEs), so the composite gain across subjects is a simple average gain where each growth measure is weighted according to the proportion of students linked to that gain. For our sample teacher, the total number of Full-Year Equivalent (FYE) students affiliated with growth measures from the gain model is 65 + 70, or 135 students. The EOG Reading grade 7 growth measure would be weighted at 65/135, and the EOG Math grade 7 growth measure would be weighted at 70/135.

The calculation of the composite gain across subjects based on the gain model is as follows:

$$\text{Composite Gain} = \frac{65}{135} \text{Read}_7 + \frac{70}{135} \text{Math}_7 = \left(\frac{65}{135}\right)(-0.30) + \left(\frac{70}{135}\right)(3.80) = 1.83 \quad (24)$$

5.1.2.2 Composite Standard Error Across Subjects

As discussed in [Section 4](#), the standard error is a measure of the statistical certainty in the growth measure that indicates whether an estimate is decidedly above or below expected growth. Standard errors can, and should, also be provided for the composite gains that have been calculated from a teacher's value-added gain estimate.

As background, statistical formulas are often more conveniently expressed as variances (see equation 6, for example), and this is the square of the standard error. Standard errors of composites can be calculated using variations of the general formula shown below. To maintain the generality of the formula, the individual estimates in the formula (think of them as value-added gains) are simply called X , Y , and Z . If there were more than or fewer than three estimates, the formula would change accordingly. As OST composites use proportional weighting according to the number of students linked to each value-added gain, each estimate is multiplied by a different weight: a , b , or c .

$$\begin{aligned} \text{Var}(aX + bY + cZ) &= a^2\text{Var}(X) + b^2\text{Var}(Y) + c^2\text{Var}(Z) \\ &+ 2ab \text{Cov}(X, Y) + 2ac \text{Cov}(X, Z) + 2bc \text{Cov}(Y, Z) \end{aligned} \quad (25)$$

Covariance, denoted by Cov , is a measure of the relationship between two variables. It is a function of a more familiar measure of relationship, the correlation coefficient. Specifically, the term $\text{Cov}(X, Y)$ is calculated as follows:

$$\text{Cov}(X, Y) = \text{Correlation}(X, Y)\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)} \quad (26)$$

The value of the correlation ranges from -1 to +1, and these values have the following meanings.

- A value of zero indicates no relationship.
- A positive value indicates a positive relationship, or Y tends to be larger when X is larger.
- A negative value indicates a negative relationship, or Y tends to be smaller when X is larger.

Two variables that are unrelated have a correlation and covariance of zero. Such variables are said to be statistically independent. If the X and Y values have a positive relationship, then the covariance will also be positive. As a general rule, two value-added gain estimates are statistically independent if they are based on completely different sets of students. For our sample teacher's composite gain, the relationship will generally be positive, and this means that the gain-based composite standard error is larger than it would be assuming independence.

For the sample teacher, it cannot be assumed that the gains in the composite are independent because it is likely that some of the same students are represented in different value-added gains, such as EOG Math grade 7 in Year 1 and EOG Reading grade 7 in Year 1.

However, to demonstrate the impact of the covariance terms on the standard error, it is useful to calculate the standard error using (inappropriately) the assumption of independence and to compare it to the standard error calculated assuming (inappropriately) an extreme correlation of +1. Using the gain-based FYE weightings and standard errors reported in Table 3 and assuming total independence, the standard error would then be as follows:

$$\begin{aligned} \text{Composite Standard Error} &= \sqrt{\left(\frac{65}{135}\right)^2 (SE_{Read_7})^2 + \left(\frac{70}{135}\right)^2 (SE_{Math_7})^2} \\ &= \sqrt{\left(\frac{65}{135}\right)^2 (1.20)^2 + \left(\frac{70}{135}\right)^2 (1.50)^2} = 0.97 \end{aligned} \tag{27}$$

Assuming a correlation of +1, the standard error would then be as follows:

$$\begin{aligned} &\text{Comp Standard Error} \\ &= \sqrt{\left(\frac{65}{135}\right)^2 (SE_{Read_7})^2 + \left(\frac{70}{135}\right)^2 (SE_{Math_7})^2 + 2\left(\frac{65}{135}\right)\left(\frac{70}{135}\right)(SE_{Read_7})(SE_{Math_7})} \\ &= \sqrt{\left(\frac{65}{135}\right)^2 (1.20)^2 + \left(\frac{70}{135}\right)^2 (1.50)^2 + 2\left(\frac{65}{135}\right)\left(\frac{70}{135}\right)(1.20)(1.50)} = 1.36 \end{aligned} \tag{28}$$

The actual standard error will fall somewhere between the two extreme values of 0.97 and 1.36 with the specific value depending on the values of the correlations between pairs of value-added gains. The magnitude of each correlation depends on the extent to which the same students are in both estimates for any two subject, grade, and year estimates.

For example, if the Year 1 EOG Math grade 7 and 2019 EOG Reading grade 7 classes had no students in common, then their correlation would be zero. On the other hand, if the Year 1 EOG Math grade 7 and Year 1 EOG Reading grade 7 classes contained many of the same students, there would be a positive correlation. However, even if those two classes had exactly the same students, the correlation would likely be considerably less than +1. Correlations of gains across years might be positive or slightly negative as the same student's score can be used in multiple gains. The actual correlations and covariances themselves are obtained as part of the modeling process using equation (10) from [Section 2.2.4.1](#). It would be impossible to obtain them outside of the modeling process. This process uses all the information about which students are in which subject, grade, and year for each teacher.

Although this approach uses a more sophisticated technique, it more accurately captures the potential relationships among teacher estimates and student scores. This will lead to the appropriate standard error that will typically be between these two extremes, which are 0.97 and 1.36 in this example. In general, the standard error of the composite gain will vary depending on the standard errors of the value-added gains and the correlations between pairs of value-added gains. The standard errors of the individual value-added gains will depend on the quantity and quality of the data that went into the gain, such as the number of students and the amount of missing data all those students have will contribute to the magnitude of the standard error.

5.1.2.3 Composite Index Across Subjects

The final step is to calculate the composite index based on the gain model, which is the composite gain divided by the composite standard error. The composite index for the sample teacher is 1.83 divided by a number between 0.97 and 1.36. The actual gain-based standard error is determined using all the information described above, which includes information beyond just our one sample teacher. For simplicity's sake, let's assume that the actual standard error in this example was 1.15, and the index for this teacher would be calculated as follows:

$$\text{Composite Index} = \frac{\text{Composite Gain}}{\text{Composite Standard Error}} = \frac{1.83}{1.15} = 1.59 \quad (29)$$

Although some of the values in the example were rounded for display purposes, the actual rounding or truncating only occurs after all the measures have been combined as described in [Section 4.4](#).

5.1.3 Technical Description of the Composite Index Based on Predictive Model Measures

For our sample teacher (and for the majority of teachers who receive reporting from the predictive model in North Carolina), there is only one available value-added measure from the predictive model. This means that the reported value-added index for that subject will be the same that is calculated for the predictive-based composite index. For the sample teacher, only a Math I growth measure is available for Year 1.

$$\text{Composite Index} = \frac{\text{Math I Growth Measure}}{\text{Math I Standard Error}} = \frac{11.75}{6.20} = 1.90 \quad (30)$$

However, should a teacher have more than one value-added measure based on the predictive model, then the composite index would be calculated by first calculating index values for each subject and then combining those weighting by the effective number of students. The standard error of this combined index must assume independence since the measures from the predictive model are done in separate models for each year and subject.

5.1.4 Technical Description of the Combined Composite Index Across Subjects Based on the Gain and Predictive Models

The two composite indices from the gain and predictive models are weighted according to the number of students linked to each model to determine the combined composite index.

Our sample teacher has 155 students of which 135 are linked to the gain model and 20 to the predictive model. The combined composite index would be calculated as follows using these weightings, the gain-based composite index across subjects, and the predictive-based index across subjects:

$$\text{Unadjusted Combined Comp Index} = \left(\frac{135}{155}\right)(1.59) + \left(\frac{20}{155}\right)(1.90) = 1.62 \quad (31)$$

This combined index is not an actual index itself until it is adjusted to accommodate for the fact that it is based on multiple pieces of evidence together. An index, by definition, has a standard error of 1, but this unadjusted value (1.62) does not have a standard error of 1. The next step is to calculate the new standard error and divide the combined composite index found above by it. This new, adjusted composite index will be the final index with a standard error of 1. The standard error can be found given the standard formula above and the fact that each index has a standard error of 1. Independence is assumed since these are done outside of the models. In this example, the standard error would be as follows:

$$\text{Final Combined Comp SE} = \sqrt{\left(\frac{135}{155}\right)^2 (1)^2 + \left(\frac{20}{155}\right)^2 (1)^2} = 0.88 \quad (32)$$

Therefore, the final combined composite index value is 1.62 divided by 0.88, or 1.85. This is the value in the teacher's SGM report.

5.2 District and School Composites

The composites calculated for schools are done in the exact same way as teachers described in the section above.

6 Input Data Used in the North Carolina Growth Model

6.1 Assessment Data Used in North Carolina

For the analysis and reporting based on the 2022-23 school year, EVAAS receives the following assessments for use in the growth and/or projection models:

- BOG Reading in grade 3
- EOG Mathematics in grades 3–8
- EOG Reading in grades 3–8
- EOG Science in grades 5 and 8
- EOC Biology, English II, Math 1 and Math 3
- Early Literacy assessments in K–2 based on mCLASS
- Career and Technical Education (CTEs) assessments in various subjects
- ACT assessments in English, Math, Reading, and Science
- PSAT assessments in Evidence-Based Reading and Writing and Math
- SAT assessments in Evidence-Based Reading and Writing and Math
- AP assessments in various subjects

The state EOG tests are administered in the spring semester whereas EOCs and CTEs are typically given in the fall and spring semesters with the occasional summer administration. BOG Reading for grade 3 is administered in the fall semester. The Early Literacy assessments are administered three times throughout the year.

Assessment files provide the following data for each student score:

- Scale score
- Test taken
- Tested grade
- Tested semester
- District number
- School number
- Membership
- Accountability growth membership
- Partial enrollment
- First Year English Learner (EL)
- First Year Math 1

Some of this information, such as performance levels, is not relevant to the ACT, PSAT, or SAT tests. At times, pre-test scores are provided for CTE assessments.

6.2 Student Information

Student information is used in creating the web application to assist educators analyze the data to inform practice and assist all students with academic growth. SAS receives this information in the form of various socioeconomic, demographic, and programmatic identifiers provided by NCDPI. Currently, these categories are as follows:

- Academically or Intellectually Gifted (Y, N)
- Gender (M, F)
- English Learners (EL) (Y, N)
- Economically Disadvantaged Students (Y, N)
- Students with Disabilities (Y, N)
- Homeless (Y, N)
- Race
 - American Indian/Alaskan Native
 - Asian/Pacific Islander
 - Black (not Hispanic)
 - Hispanic
 - Two or More Races
 - White (not Hispanic)

6.3 Teacher Information

In order to provide accurate and verified student-teacher linkages in the teacher growth models, North Carolina educators are given the opportunity to complete roster verification. This process enables teachers to confirm their class rosters for students in a particular subject, grade, and year, and it captures scenarios where multiple teachers have instructional responsibility for students. Administrators also verify the linkages as an additional check. Roster verification, therefore, increases the reliability and accuracy of teacher-level analyses.

NCDPI sends SAS teacher identification data and student-teacher linkages from PowerSchool. The student-teacher linkage files include the following information:

- Teacher identification
 - Teacher Name
 - Teacher Unique ID
- Student linking information
 - Student Last Name
 - Student First name
 - Unique Student ID (USID)
- Course information linked to a tested subject via a course to subject mapping provided by DPI
- District and School information (numbers)
- Percentage of instructional responsibility derived from enrollment information provided by DPI (i.e., the date the student enrolled and the date the student left the course).

7 Business Rules

7.1 Assessment Verification for Use in Growth Models

To be used appropriately in any growth models, the scales of these assessments must meet three criteria:

1. **There is sufficient stretch in the scales** to ensure progress can be measured for both low-achieving students as well as high-achieving students. A floor or ceiling in the scales could disadvantage educators serving either low-achieving or high-achieving students.
2. **The test is highly related to the academic standards** so that it is possible to measure progress with the assessment in that subject, grade, and year.
3. **The scales are sufficiently reliable from one year to the next.** This criterion typically is met when there are a sufficient number of items per subject, grade, and year. This will be monitored each subsequent year that the test is given.

These criteria are checked annually for each assessment prior to use in any growth model, and North Carolina's current standardized assessments meet them. These criteria are explained in more detail below.

7.1.1 Stretch

Stretch indicates whether the scaling of the assessment permits student growth to be measured for both very low- or very high-achieving students. A test "ceiling" or "floor" inhibits the ability to assess students' growth for students who would have otherwise scored higher or lower than the test allowed. It is also important that there are enough test scores at the high or low end of achievement, so that measurable differences can be observed.

Stretch can be determined by the percentage of students who score near the minimum or the maximum level for each assessment. If a much larger percentage of students scored at the maximum in one grade than in the prior grade, then it might seem that these students had negative growth at the very top of the scale when it is likely due to the artificial ceiling of the assessment. Percentages for all North Carolina assessments are well below acceptable values, meaning that these assessments have adequate stretch to measure value-added even in situations where the group of students are very high or low achieving.

7.1.2 Relevance

Relevance indicates whether the test is sufficiently aligned with the curriculum. The requirement that tested material correlates with standards will be met if the assessments are designed to assess what students are expected to know and be able to do at each grade level. This is how state tests are designed and is monitored by the NCDPI and their psychometricians.

7.1.3 Reliability

Reliability can be viewed in a few different ways for assessments. Psychometricians view reliability as the idea that a student would receive similar scores if the assessment was taken multiple times. The type of reliability is important for most any use of standardized assessments.

7.2 Pre-Analytic Processing

7.2.1 Missing Grade

In North Carolina, the grade used in the analyses and reporting is the tested grade, not the enrolled grade. If a grade is missing on an early grade or end-of-grade test record, then that record will be excluded from all analyses. The grade is required to include a student's score in the appropriate part of the models and to convert the student's score into the appropriate NCE in the gain-based model.

7.2.2 Duplicate (Same) Scores

If a student has a duplicate score for a particular subject and tested grade in a given testing period in a given school, then the extra score will be excluded from the analysis.

7.2.3 Students with Missing Districts or Schools for Some Scores but Not Others

If a student has a score with a missing district or school for a particular subject and grade in a given testing period, then the duplicate score that has a district and/or school will be included over the score that has the missing data.

7.2.4 Students with Multiple (Different) Scores in the Same Testing Administration

If there are multiple records within a year/semester, the analysis only includes the earliest test administration date. If a student has multiple scores in the same period for a particular subject and grade and the test scores are not the same, then those scores will be excluded from the analysis.

Priority is given to assessment records used in prior analyses since these historical records passed through and are approved for use in the statewide value-added analysis.

If duplicate scores for a particular subject and tested grade in a given testing period are at different accountable schools, then both scores will be excluded from the analysis.

The highest composite combination of SAT subjects is used for SAT value-added and student college readiness projections.

Note that if multiple scores are received for grade 3 Reading or Math across years, only the most recent score is used.

7.2.5 Students with Multiple Grade Levels in the Same Subject in the Same Year

A student should not have different tested grade levels in the same subject in the same year. If that is the case, then the student's records are checked to see whether the data for two separate students were inadvertently combined. If this is the case, then the student data are adjusted so that each unique student is associated with only the appropriate scores. If the scores appear to all be associated with a single unique student, then scores that appear inconsistent are excluded from the analysis. For the K-2 Early Literacy Assessments, if a student's tested grade changes within a single year, then the analysis excludes both records in that year.

7.2.6 Students with Records That Have Unexpected Grade Level Changes

If a student skips more than one grade level (e.g., moves from sixth in 2018 to ninth in 2019) or is moved back by one grade or more (i.e., moves from fourth in 2018 to third in 2019) in the same subject, then the student's records are examined to determine whether two separate students were inadvertently

combined. If this is the case, then the student data is adjusted so that each unique student is associated with only the appropriate scores. These scores are removed from the analysis if it is the same student. Per DPI's decision, the analysis does not remove students with scores that appear to be associated with inconsistent grades. The analysis leaves students in the analysis at the tested grade that EVAAS receives from DPI.

7.2.7 Students with Historical EOC Records that Occur Too Early or Late Given the Student's Testing History

Student records are excluded for students who have historical EOC records that occur too early or late given the student testing history. The last observed grade level for each student on grade-level assessments (i.e., EOGs) is used to determine the student's associated cohort. Based on this cohort, the analysis determines a likely expected grade for the student when they take their EOC assessments. If the likely expected grade for the EOGs is 4 or less for Math 1, 5 or less for Biology or English II, or 15 or greater for any EOC including Math 3, then the record excluded from analyses because the predictors available for those EOC subjects are too far removed from the EOC administration or potentially connected to a student with questionable identifying information. Note that grade 15 is the equivalent of a student who was expected to graduate over three years ago.

7.2.8 Students with Records at Multiple Schools in the Same Test Period

If a student is tested at two different schools in a given testing period, then the student's records are examined to determine whether two separate students were inadvertently combined. If this is the case, then the student data is adjusted so that each unique student is associated with only the appropriate scores. When students have valid scores at multiple schools in different subjects, all valid scores are used at the appropriate school.

7.2.9 Outliers

Student assessment scores are checked each year to determine whether they are outliers in context with all the other scores in a reference group of scores from the individual student. These reference scores are weighted differently depending on proximity in time to the score in question. Scores are checked for outliers using related subjects as the reference group. For example, when searching for outliers for Math test scores, all EOG and EOC Math subjects are examined simultaneously, and any scores that appear inconsistent, given the other scores for the student, are flagged. Outlier identification for college readiness assessments use all available college readiness data alongside state assessments in the respective subject area (e.g., Math subjects with EOC, EOG, and PSAT tests might be used to identify outliers with SAT or ACT). Furthermore, K-2 Early Literacy data are used solely for outlier identification with Early Literacy assessments. Lastly, CTE assessments do not undergo outlier identification due to the various test taking patterns inherent with CTE and the fact that these assessments have less uniformity in administration across the state than other statewide assessments.

Scores are flagged in a conservative way to avoid excluding any student scores that should not be excluded. Scores can be flagged as either high or low outliers. Once an outlier is discovered, that outlier will not be used in the analysis, but it will be displayed on the student testing history on the EVAAS web application.

This process is part of a data quality procedure to ensure that no scores are used if they were, in fact, errors in the data, and the approach for flagging a student score as an outlier is fairly conservative.

Considerations included in outlier detection are:

- Is the score in the tails of the distribution of scores? Is the score very high or low achieving?
- Is the score “significantly different” from the other scores as indicated by a statistical analysis that compares each score to the other scores?
- Is the score also “practically different” from the other scores? Statistical significance can sometimes be associated with numerical differences that are too small to be meaningful.
- Are there enough scores to make a meaningful decision?

To decide whether student scores are considered outliers, all student scores are first converted into a standardized normal Z-score. Then each individual score is compared to the weighted combination of all the reference scores described above. The difference of these two scores will provide a t-value of each comparison. Using this t-value, the growth models can flag individual scores as outliers.

There are different business rules for the low outliers and the high outliers, and this approach is more conservative when removing a very high-achieving score.

For low-end outliers, the rules are:

- The percentile of the score must be below 50.
- The t-value must be below -3.5 for EOGs and EOCs when determining the difference between the score in question and the weighted combination of reference scores (otherwise known as the comparison score). In other words, the score in question must be at least 3.5 standard deviations below the comparison score. For other assessments, the t-value must be below -4.0
- The percentile of the comparison score must be above a certain value. This value depends on the position of the individual score in question but will range from 10 to 90 with the ranges of the individual percentile score.

For high-end outliers, the rules are:

- The percentile of the score must be above 50.
- The t-value must be above 4.5 for EOGs and EOCs when determining the difference between the score in question and the reference group of scores. In other words, the score in question must be at least 4.5 standard deviations above the comparison score. For other assessments, the t-value must be above 5.0.
- The percentile of the comparison score must be below a certain value. This value depends on the position of the individual score in question but will need to be at least 30 to 50 percentiles below the individual percentile score.
- There must be at least three scores in the comparison score average.

7.2.10 Linking Records Over Time

Each year, EVAAS receives data files that include student assessment data and file formats. These data are checked each year prior to incorporation into a longitudinal database that links students over time. Student test data and demographic data are checked for consistency year to year to ensure that the appropriate data are assigned to each student. Student records are matched over time using all data provided by the state, and teacher records are matched over time using the Unique ID and teacher’s name.

7.3 Growth Models

7.3.1 Students Included in the Analysis

As described in Pre-Analytic Processing, student scores might be excluded due to the business rules, such as outlier scores.

For the gain model, all students are included in these analyses if they have assessment scores that can be used. The gain model uses all available EOG Math and Reading results for each student. Because this model follows students from one grade to the next and measures growth as the change in achievement from one grade to the next, the gain model assumes typical grade patterns for students. Students with non-traditional patterns, such as those who have been retained in a grade or skipped a grade, are treated as separate students in the model. In other words, these students are still included in the gain model, but the students are treated as separate students in different cohorts when these non-traditional patterns occur. This process occurs separately by subject since some students can be accelerated in one subject and not in another. Students are excluded from the gain model if the student is flagged as a First Year EL student or if the student does not meet partial enrollment membership.

For the predictive and projection models, a student must have at least three valid predictor scores that can be used in the analysis, all of which cannot be deemed outliers. (See [Section 7.2.9](#) on Outliers.) These scores can be from any year, subject, and grade that are used in the analysis. In other words, the student's expected score can incorporate other subjects beyond the subject of the assessment being used to measure growth. The required three predictor scores are needed to sufficiently dampen the error of measurement in the tests to provide a reliable measure. If a student does not meet the three-score minimum, then that student is excluded from the analyses. It is important to note that not all students have to have the same three prior test scores; they only have to have some subset of three that were used in the analysis. Unlike the gain model, students with non-traditional grade patterns are included in the predictive model as one student. Since the predictive model does not determine growth based on consecutive grade movement on tests, students do not need to stay in one cohort from one year to the next. That said, if a student is retained and retakes the same test, then that prior score on the same test will not be used as a predictor for the same test as a response in the predictive model. This is mainly due to the fact that very few students used in the models have a prior score on the same test that could be used as a predictor. In fact, in the predictive model, it is typically the case that a prior test is only considered a possible predictor when at least 50% of the students used in that model have those prior test scores. Students are excluded from the predictive model if the student is flagged as a First Year EL student or if the student does not meet partial enrollment membership for EOG Science, EOC, and CTE assessments. There are no membership rules used to include or exclude students in the SAT, PSAT, and ACT analyses.

7.3.2 Minimum Number of Students to Receive a Report

The growth models require a minimum number of students in the analysis in order for districts, schools, and teachers to receive a growth report. This is to ensure reliable results.

7.3.2.1 District and School Model

For the gain model, the minimum student count to report an estimated average NCE *score* (i.e., either entering or exiting achievement) is six students in a specific subject, grade, and year. To report an estimated NCE *gain* in a specific subject, grade, and year, there are additional requirements:

- There must be at least six students who are associated with the school or district in the subject, grade, and year.
- Of those students who are associated with the school or district in the current year and grade, there must be at least six students in each subject, grade, and year in order for that subject, grade, and year to be used in the gain calculation.
- There is at least one student at the school or district who has a “simple gain,” which is based on a valid test score in the current year and grade as well as the prior year and grade in the same subject. However, due to the rule above, it is typically the case that at least six students have a “simple gain.” In some cases where students only have a Math or Reading score in the current year or previous year, this value dips below six.
- For any district or school growth measures based on specific student groups, the same requirements described above apply for the students in that specific student group.

For example, to report an estimated NCE gain for school A in EOG Math grade 5 for this year, there must be the following requirements:

- There must be at least six fifth-grade students with an EOG Math grade 5 score at school A for this year.
- Of the fifth-grade students at school A this year *in all subjects, not just Math*, there must be at least six students with an EOG Math grade 4 score from last year.
- At least one of the fifth-grade students at school A this year must have an EOG Math grade 5 score from this year *and* an EOG Math grade 4 score from last year.

For the predictive model, the minimum student count to receive a growth measure is 10 students in a specific subject, grade, and year. These students must have the required three prior test scores needed to receive an expected score in that subject, grade, and year.

For any district or school growth measures based on specific student groups, the same requirements described above apply for the students in that specific student group. Note that while subgroup reporting requires 10 students to be included in the EVAAS web reporting, it requires 30 students to be included in the state’s accountability model for student groups.

7.3.2.2 Teacher Model

The teacher gain *model* includes teachers who are linked to at least six students with a valid test score in the same subject, grade, and year. This requirement does not consider the percentage of instructional time that the teacher spends with each student in a specific subject and grade.

The teacher predictive *model* includes teachers who are linked to at least 10 students with a valid test score in the same subject/grade or course within a year. This requirement does not consider the percentage of instructional time the teacher spends with each student in a specific subject and grade.

For both the gain and predictive models, to receive a Teacher *report* in a particular year, subject, and grade, there is an additional requirement. A teacher must have at least six Full Time Equivalent (FTE) students in a specific subject, grade, and year. The teacher's number of FTE students is based on the number of students linked to that teacher and the percentage of instructional time the teacher has for each student. For example, if a teacher taught 10 students for 50% of their instructional time, then the teacher's FTE number of students would be five, and the teacher would not receive a teacher growth report. If another teacher taught 12 students for 50% of their instructional time, then that teacher would have six FTE students and would receive a teacher growth report. The instructional time attribution is obtained from the linkage roster verification process that is used in North Carolina.

The teacher gain model has an additional requirement. The teacher must be linked to at least five students, and one of these five students must have a "gain," meaning the same subject prior test score must come from the immediate prior year and prior grade. Note that if a student repeats a grade, then the prior test data would not apply as the student has started a new cohort.

7.3.2.3 Student Groups

The student groups model has the same requirements as the district/school models for the minimum number of students.

7.3.3 Student-Teacher Linkages

Student-teacher linkages are connected to assessment data based on the subject and identification information described in [Section 6.3](#). The model will make adjustments to linkages if a student is claimed by teachers at a total percentage higher than 100% in an individual year, subject, and grade. If over-claiming happens, then the individual teacher's weight is divided by the total sum of all weights to redistribute the attribution of the student's test scores across teachers. Underclaimed linkages for students are not adjusted because a student can be claimed less than 100% for various reasons (such as a student who lives out of state for part of the year).