

Article

Statistical modelling and analysis of NCEA and New Zealand Scholarship assessment data

Michael Johnston^{1*} and David Lillis²

¹Faculty of Education, Victoria University of Wellington, PO Box 600, Wellington and

²New Zealand Qualifications Authority, PO Box 160, Wellington 6140

New Zealand's main qualification system for senior secondary school comprises the three levels of the National Certificate of Educational Achievement (NCEA). These qualifications were introduced progressively, Level 1 first becoming available in 2002, Level 2 in 2003, and Level 3 in 2004. Additionally, the present system for awarding New Zealand Scholarship was first implemented in 2005. The NCEA system has several features that are quite unique, and that afford schools the opportunity to develop their own assessment programmes for a wide variety of courses in traditional, emerging, and cross-disciplinary subject areas. Those features of the NCEA system that afford this flexibility also present challenges from the psychometric point of view. In this paper, we describe a range of statistical modelling and analyses undertaken by the New Zealand Qualifications Authority (NZQA) to meet these challenges.

The NCEA system is criterion-referenced rather than norm-referenced. This means that assessment results depend on the performance of candidates against set criteria, rather than being determined relative to the performance of other candidates.

Criterion-referencing is not unique to NCEA. Indeed, a movement from norm-referenced to criterion-referenced assessment is evident in many assessment systems around the world (e.g. Australia and the United States). Criterion-referenced assessment results and qualifications arguably carry greater inherent meaning than those based on norm-referencing because, if a candidate meets the criterion for a particular result, it is possible, within the limits of the reliability and validity of the assessment, to certify him or her as competent in the specific skill or knowledge associated with that criterion. Under a norm-referenced system, the only information that can be inferred validly from a

candidate's result are the percentages of other candidates who demonstrated higher or lower performance.

From a psychometric perspective, running a high-quality criterion-referenced system is more challenging than running a norm-referenced system. Under the latter, all that is required is an accurate rank-order of the candidates, with normative scaling used to allocate final results on the basis of that rank order. Differences in the difficulty of an assessment (for example, a formal examination) from year to year do not affect outcomes unless these would result in a different rank ordering of candidates.

Under a criterion-referenced system, however, the standard of performance commensurate with the criteria must be maintained over time. Under any assessment system the connection between candidates' performance in an assessment and the final results must entail expert judgement, and cannot be established on a purely statistical basis. However, in a large-scale criterion-referenced system such as NCEA, professional judgement requires a great deal of statistical and psychometric support if criteria are to be applied consistently across different assessors and over time.

Perhaps the most unique aspect of NCEA is its decomposition of assessment into units known as 'standards'. Whereas, under most secondary assessment systems internationally, candidates receive a single result for each subject they have studied, under NCEA candidates receive multiple results, each certifying specific skills and knowledge. For example, there is a trigonometry standard, called *Solve right-angled triangle problems*, and another pertaining to English-language literacy called *Read and understand unfamiliar texts*.

*Correspondence: Michael.Johnston@vuw.ac.nz



Michael Johnston has recently commenced as a senior lecturer in the School of Educational Policy and Implementation at Victoria University. He was previously a senior statistician at the New Zealand Qualifications Authority, where he conducted research, analysis and evidence-based policy development for a range of reforms to assessment systems for NCEA and New Zealand Scholarship. Dr Johnston qualified for his PhD at the University of Melbourne. He has extensive experience in experimental psychology and other quantitative research in social science and education. He is a member of the New Zealand Assessment Academy and of the Technical Overview Group (Assessment), an independent committee of academics providing technical advice to NZQA.

David Lillis is a senior statistician with the New Zealand Qualifications Authority (NZQA). He holds a PhD from Curtin University in Western Australia. At NZQA he conducts a wide range of data analysis, including the analysis of NCEA and New Zealand Scholarship results. In particular, he writes software in the R language for Item Response Theory as one approach to ensuring the high quality of secondary examinations. Dr Lillis is a past president of the New Zealand Association of Scientists.



It is this aspect of NCEA that affords its great flexibility, because schools can choose standards that best reflect the content of their courses, and can assess cross-disciplinary courses by selecting relevant standards from more than one subject area. Nonetheless, maintaining consistency of assessment judgements over the approximately 700 standards that are derived from the New Zealand curriculum presents a difficult psychometric problem. In part, this is because there are so many standards, but mainly it is because the assessment for each standard is necessarily of shorter duration and entails a smaller volume of work than would be the case if assessment were conducted at the level of the subject. The difficulty that this situation presents is one of maintaining assessment reliability - shorter and smaller-volume assessments tend to have poorer reliability than longer or larger volume assessments (assuming similar assessment quality).

In this paper we describe a number of statistical processes that assist NZQA to meet the challenges posed by the design of NCEA in relation to external assessment; that is, assessment procedures designed and administered by NZQA, a large majority of which are time-limited examinations. Internal assessments, those designed and conducted in schools and moderated by NZQA, also comprise a very important component of NCEA, and NZQA does have procedures for monitoring the reliability of teachers' internal assessment judgements. However, discussion of these procedure is beyond the scope of the present paper.

The processes we discuss here are as follows: the development and use of Profiles of Expected Performance (PEPs), used as a guide to maintain standards during the marking of external assessments; a set of post-hoc analyses of NCEA examination results, carried out annually following each external assessment round in order to assess the performance of examination items and papers; and statistical procedures used to assist in the allocation of results for New Zealand Scholarship assessments, as well as analysis of the quality of these examinations.

The analyses described here are used to inform, rather than replace, expert judgement. Collectively, these procedures provide assessment practitioners with support for their professional judgement, and with information that enables them to maintain and improve their consistency in applying the various assessment criteria of each standard.

Profiles of expected performance

In the early years of NCEA it was found that, for many externally assessed standards, the proportions of candidates receiving each grade fluctuated from year to year. Given that the system was very new, some variations were to be expected. However, the size of the variation was, in many cases, large even in light of the circumstances. It soon became evident that some form of statistical support for professional judgement was required to maintain consistency in the application of the standards over time.

Profiles of Expected Performance (PEPs) were introduced in 2005 to address the problem of variations from the expected results distributions from one year to the next. The PEP gives a percentage range into which each grade – *Not Achieved* (N), *Achieved* (A), *Merit* (M) and *Excellence* (E) – is expected to fall. For example, we might expect that in a given standard 20–32% of candidates will earn an *Achieved* grade, or that 6–10% will receive *Excellence*. Figure 1 shows the 2010 PEP bands for the Level 3 Calculus standard 90636 (Integrate functions and use integrals to solve problems).

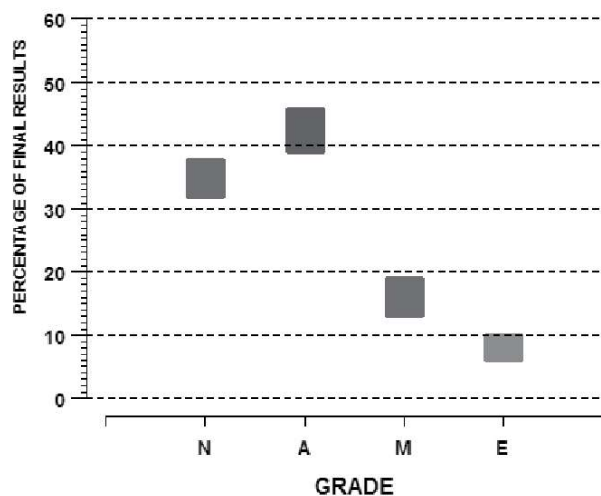


Figure 1. Profiles of Expected Performance for the Level 3 Standard: Integrate functions and use integrals to solve problems (2010).

It is not the intention of the PEP process to manipulate results to fit a pre-determined distribution. Rather, the expected statistical stability of distributions of large numbers of results is used to identify discrepancies that might signal a variation in the standard of performance required for particular grades. It is quite permissible for actual results to fall outside PEP ranges. However, when this occurs, there must be a defensible explanation for the discrepancy that does not entail any implicit change in the performance criterion.

If, during marking, it appears that any of the grades will fall outside the expected range for a particular standard, a discussion is held between NZQA and the leader of the marking panel to discuss reasons for the difference. If there is a legitimate reason (for example, that the characteristics of the cohort have changed in some way, or that there has been an overall improvement or deterioration in performance), then the distribution stands unchanged. If, on the other hand, the reason does not appear to be legitimate, then the marking schedule may be revised. For example, an easier examination than those of previous years is not an acceptable reason for result falling outside PEPs; notwithstanding the difficulty of an examination, candidates must meet the same standard each year in order to receive a particular grade.

A PEP is generated for each grade in each externally-assessed standard in which at least 300 candidates have entered. Below this number, the statistical stability of distributions of results is insufficient to justify the development of a PEP. All PEPs are set prior to each year's examination round, taking into account the history of results for the standard, as well as statistical estimates of the distribution expected on the basis of the previous year's candidature across other standards.

PEPs for standards with large cohorts are set with tighter confidence bands than those with smaller cohorts. Small cohorts lead to lower stability than large cohorts. A substantial change in cohort size from the previous year may also justify setting a larger confidence band, because usually it is not possible to predict in advance the characteristics of the larger cohort.

Draft PEPs are set initially on the basis of the history of results for the standard, as well as professional knowledge of the subject area and candidature. Usually the PEPs for a standard

will be the same or very similar from year to year. Following the development of the draft PEPs, other statistical information is taken into account, perhaps prompting a revision of the draft. This statistical information includes analysis of the difficulty of the standard and the overall ability of the cohort, based on the previous year's results.

Measurement of the difficulty of a standard T_i involves comparing the relative performance of candidates undertaking T_i with their performance on each other standard $T_{j1}, T_{j2}, \dots, T_{jn}$ that has an overlapping cohort (i.e. a set of candidates who undertake both assessments) with T_i . Figure 2 gives a diagram of this situation: there is a target standard, T_i , and two other standards with overlapping cohorts: T_{j1} and T_{j2} . (In a real world example there would be many overlapping standards.) The overlapping cohorts for the pairs T_i, T_{j1} , and T_i, T_{j2} are labelled c_{ij1} and c_{ij2} respectively.

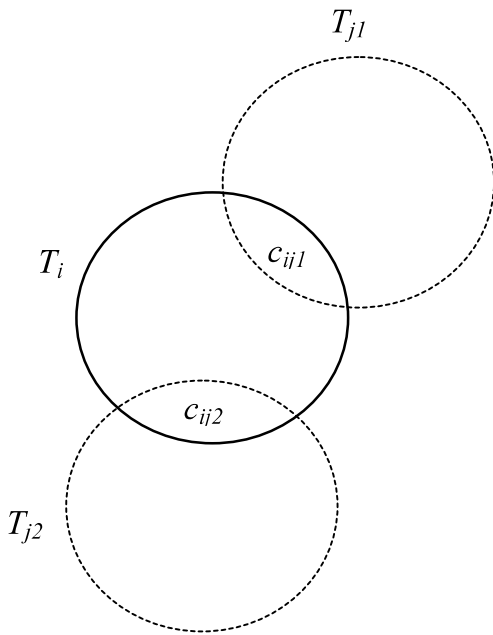


Figure 2. Diagram depicting overlapping cohorts for a target standard, T_i , and two other standards, T_{j1} and T_{j2} . The overlapping cohorts are designated c_{ij1} and c_{ij2} respectively.

Equation 1 provides a formal method for determining the difficulty of a target standard relative to other standards with overlapping cohorts. The difficulty of T_i compared to other standards can be estimated by calculating a mean difference in the rate of success for the cohort c_{ij1} on T_i and the rate of success for the same cohort on each other standard T_j .

The average differences in rates of success are in fact weighted averages, in which the magnitudes of the weights are determined by the relative sizes of the overlapping cohorts and by the correlation in performance between the target standard and each overlapping standard. Weighting by the size of the overlap places greater emphasis on comparisons involving standards with larger common cohorts, because larger overlaps result in more reliable comparisons.

The correlation in rate of success measures the extent to which performance in a pair of standards draws upon similar knowledge, skills, or cognitive functions. Clearly, if perform-

ance in two standards is uncorrelated (i.e. if the value of the correlation coefficient is zero), then the question of their relative difficulty does not arise. On the other hand, if performance in two standards were completely correlated (i.e. the value of the correlation coefficient were unity), then performance on one would be completely predictable from performance on the other, and they would be fully comparable in difficulty. In practice, correlations are never perfect, and although the theoretical minimum correlation is negative one (a negative correlation indicating an inverse relationship in performance), correlations in performance on pairs of standards as low as zero are very rarely, if ever, observed.

Equation 1 gives a mathematical expression that is used to calculate the relative difficulty of a standard using information on candidate performance across all standards held on NZQA's results databases.

$$D_i = \frac{\sum_{j=1}^n c_{ij} \rho_{ij} [R_{ij}(j) - R_{ij}(i)]}{\sum_{j=1}^n c_{ij}}$$

Equation 1. Difficulty (D) of a standard i , where c_{ij} is the number of candidates undertaking both standard i and each other standard j , ρ_{ij} is the magnitude of the correlation (Spearman's ρ) between standard i and each other standard j , $R_{ij}(i)$ is the rate of success in standard i of the overlapping cohort, $R_{ij}(j)$ is the rate of success of the overlapping cohort in standard j , and n is the total number of standards with cohorts overlapping that of standard i .

If the success rate in standard i is high (i.e. the standard is easier than an overlapping standard j), then the success rate of the overlapping cohort in that standard, $R_{ij}(i)$, is higher than the success rate of that cohort in the overlapping standard $R_{ij}(j)$. In this case the difference $R_{ij}(j) - R_{ij}(i)$ is negative and decreases D_i slightly. Conversely, standards that are difficult relative to comparison standards increase the magnitude of D_i . The denominator is the sum of all cohort sizes and is intended to constrain the magnitude of D_i to a useful range of values.

The cohort strength uses a slightly different comparison (see Equation 2).

$$S_i = \frac{\sum_{j=1}^n c_{ij} \rho_{ij} [R_{ij}(i) - R_j(j)]}{\sum_{j=1}^n c_{ij}}$$

Equation 2. Strength (S) of a cohort in standard i , where c_{ij} is the number of candidates undertaking both standard i and each other standard j , ρ_{ij} is the magnitude of the correlation (Spearman's ρ) between standard i and each other standard j , $R_{ij}(i)$ is the rate of success in standard i of the overlapping cohort, $R_j(j)$ is the rate of success in standard j of candidates undertaking standard j but not standard i , and n is the total number of standards with cohorts overlapping that of standard i .

In this case, rather than comparing rates of success of a cohort in a target standard with rates of success in other standards, we compare the performance of the cohort undertaking both the

target standard and each comparison standard, with the cohort undertaking the comparison standard only. If the cohort of the target standard is strong, then any subset of that cohort (that subset overlapping with comparison standards) will tend to have a higher rate of success on that standard than the cohort taking the other standards only. In this case the difference in rates of success will be positive and the estimate of cohort strength will be commensurately high.

Post-hoc analysis of NCEA external assessments (examinations)

Every year NZQA undertakes a variety of statistical analysis and modelling of NCEA examination results, to contribute to continuous improvement of the quality of examination items and papers, and marking procedures. These analyses include tests of the dimensionality of the examinations and the inter-correlations of the examination items (questions) in order to determine the extent to which they measure on a single continuum of performance. Further analyses use a specialised branch of psychometric statistics, *Item Response Theory* (IRT), to determine the extent to which examination items are of appropriate difficulty and that they discriminate sufficiently between candidates of varying abilities.

For each examination, a sample of results from 700 examination scripts, or as many as are available, is analysed, focusing both on the performance of each item and on the examination as a whole. The analyses are designed to assist examiners in developing future examinations, and to develop items that measure candidates' performance consistently, both with respect to the standards and with respect to other items.

External assessments (examinations) for NCEA are designed to assess on a single dimension of performance, so that a single criterion for each grade is located on that single dimension. This is in part because there are many standards, resulting in a relatively short examination time for each standard. Some are examined in as little as 40 minutes, although from 2013 the minimum examination time for any standards will be one hour. From a purely statistical perspective, measurement on a single dimension requires that the data (candidates' item-level results) can be fitted to a single (quantitative) scale. In fact, the IRT techniques used to assess the difficulty and discrimination of each item are predicated on uni-dimensionality.

We use Principal Components Analysis, a technique first discussed by Pearson (1901), to explore the dimensionality of the external assessments as reflected in candidates' item grades. Principal Components Analysis is a widely-used dimension reduction technique in which observations of correlated variables are expressed as linear combinations of those variables, each combination constituting a principal component (or dimension).

Each principal component accounts for a proportion of the total variance in the data. The first accounts for the greatest variance, and subsequent principal components account for progressively smaller proportions. One approach to depicting principal components graphically is the scree plot (Cattell 1966). Figure 3 shows a scree plot for the item-level results for a sample of 597 scripts from the 2010 Level 1 Biology examination for standard 90168 (Describe biological ideas relating to how humans use and are affected by micro-organisms).

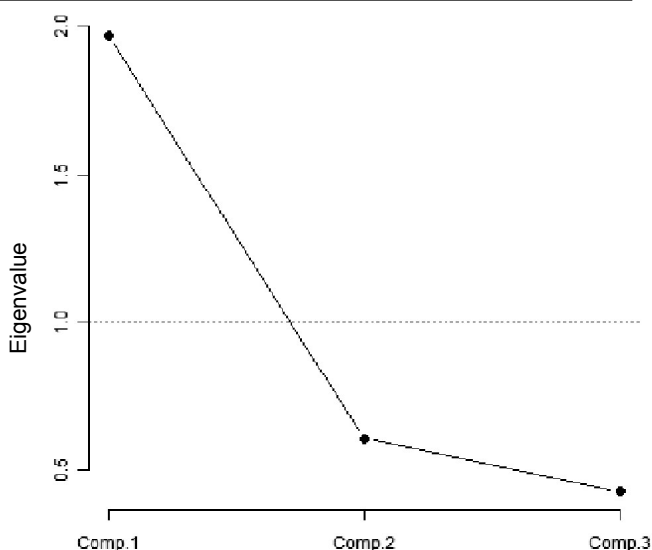


Figure 3. A scree plot showing the factor structure for the three item examination for Biology 90168 in 2010.

This particular examination comprised three items.

The vertical axis of the plot measures the magnitudes of the principal component eigenvalues. Eigenvalue magnitudes are proportional to the total variance in the data explained by each dimension. The horizontal axis of the plot displays each of the possible principal components or dimensions, one for each of the three items, arranged in order of decreasing magnitude.

How many significant dimensions (i.e. different kinds of skill or knowledge) are represented in Figure 3? One commonly-used criterion is that the eigenvalue of a significant dimension should be greater than 1. This criterion was proposed initially by Kaiser (1960), although other criteria for judging the significance of principal components have been suggested, often based on ratios of the first few eigenvalues. Because NCEA external examinations are designed to measure on a single continuum of performance, we expect only the first eigenvalue to explain a substantial fraction of the total variance.

The plot of Figure 3 suggests the presence of just one significant principal component, thus confirming the suitability of the data for the item response. We can identify those items that contribute to a particular dimension by examining the factor loadings (the correlation coefficients between the variables and principal components). Table 1 gives the factor loading of each item of the Biology examination on the first principal component.

Table 1. Item loadings on the first principal component for Biology 90168 (2010 examination round).

Item	Correlation with component 1
Q1	0.59
Q2	0.54
Q3	0.60

Loadings close to unity indicate strong relationships between the items and the components. If the examination results indicate only one dominant dimension, then most or all of the items have loaded strongly on the first principal component. Loadings above about 0.4 indicate substantial correlation with a principal component or dimension. Table 1 shows that the three items

of Biology 90168 all had moderately strong, and very similar, loadings on the first, dominant component.

For the purpose of quantitative analysis, we can treat the items on any examination that measures a single dimension as forming a distinct scale (i.e. a set of related items that measure collectively an aggregate of responses over those items). The squared factor loading gives the proportion of variance in the item results explained by a factor.

Table 2 shows two further measures of the internal consistency (or how closely related a set of item responses are when taken as a group) for the three items of the same examination for Biology 90168. These measures are the inter-item correlations and item-total correlations. They complement Principal Components Analysis in helping us to quantify the consistency of the item results and to establish the dimensionality of the examination. Both of these measures range from -1.0 to 1.0, though in practice we never encounter negative correlations between items.

Table 2. Inter-item and Item Correlations for the three items (Q1 - Q3) comprising the assessment for Biology 90168 in 2010.

Item	Q1	Q2	Q3	Total
Q1	1.00	0.43	0.58	0.60
Q2	0.43	1.00	0.44	0.49
Q3	0.58	0.44	1.00	0.60

Inter-item correlations indicate the strength of the relationships between pairs of items. Any two items that belong to the same dimension tend to exhibit strong inter-item correlation. Correlations between about 0.4 and 0.7 are considered optimal. Very high correlations (say about 0.85 or more) suggest redundancy (i.e. that we could have assessed the candidates' skills and knowledge with the same reliability using a shorter examination based on fewer items). From Table 2 we see that the correlations for the Biology standard 90168 are in this optimal range.

The item-total correlation for each item is given in the final column of Table 2. This measure is the correlation between the responses for each item and the sum of the responses for the remaining items. The item-total correlation assists in the identification of any items that are not consistent with the other items of the assessment scale. A value below 0.4 is taken as an indication that the item does not correlate well with the scale overall. In the development of psychometric tests and surveys, often such items are removed entirely. For the items of Table 2, we see that the item-total correlations of the Biology examination lie well above this threshold.

The third measure of internal consistency that we use for NCEA and New Zealand Scholarship is Cronbach's alpha (Cronbach 1951), another commonly-used measure, also ranging between -1.0 and 1.0. Cronbach's alpha can be expressed as a function of the number of test items and the average inter-correlation among the items. Cronbach's alpha tends to increase as the inter-correlations among the items increase.

The ideal range for Cronbach's alpha is from about 0.7 to about 0.85, values greater than 0.85 indicating strong homogeneity and possibly redundant items. Redundant items do not provide additional information about candidates, but simply add to the length of the assessment or test. Values substantially lower than 0.7 indicate that some items are not measuring on the same dimension as the examination as a whole.

Item Response Theory

Item Response Theory refers to a family of statistical models used to assess the quality of psychometric tests and assessments. IRT is used to inform the design, analysis and scoring of tests, questionnaires and assessment instruments, and measures abilities, attitudes and other latent traits. It is widely used internationally in the development and analysis of educational assessments.

The parameters of interest to NZQA are the difficulty of attaining a particular grade for each item, and the item discrimination, which measures how well an item discriminates between candidates of different abilities. A third parameter of interest is the ability, a measure of each candidate's performance across the entire examination (see a later section for a discussion of the ability parameter).

We use IRT to investigate the quality of our externally-assessed standards, and have developed several related approaches for conducting these analyses. Currently, we use a two-parameter graded-response model (Samejima 1969) to estimate both candidates' abilities and item parameters (discrimination and the difficulty of each assessment grade). Here, the probability of obtaining a particular grade or better (*Not Achieved*, *Achieved*, *Merit*, or *Excellence*), for a candidate of ability θ , is given by equation 3:

$$P_j(\theta) = \frac{\exp [ka(\theta - b_j)]}{1 + \exp [ka(\theta - b_j)]}$$

Equation 3. Probability of achieving a particular grade or better for a candidate of ability θ under Samejima's Graded Response Model (1969) on an item of difficulty b_j and discrimination a and where $k = -1.7$.

In equation 3 the subscript j indexes the assessment grades *Achieved* (A) or better, *Merit* (M) or better, and *Excellence* (E), θ is the calculated ability (which you can also think of as a measure of performance), P_j is the probability of achieving a particular grade or better for a candidate of ability θ , a is the fitted item discrimination, and b_j is the estimated difficulty of gaining either an A or better, M or better, or an E grade for the item. Equation 3 describes a logistic curve, and the constant k takes a value of 1.7, which scales the logistic curve such that it closely approximates a cumulative ogive. In the two-parameter model we are required to estimate the parameters a and each b_j (four parameters in total), in addition to candidates' ability parameters (one for each candidate).

Candidate ability

Ability is a multi-dimensional concept, and cannot be measured uniquely for any person. In fact, the constructs we wish to measure, such as mathematical, scientific or linguistic abilities, are actually a synthesis of many related abilities and skills. Abilities are calculated for each candidate on the basis of the entire complement of item grades. In fact, abilities estimated from IRT can provide better measures of performance than aggregates of marks or raw grade point averages, because ability estimates take explicit account of the discriminative and difficulty properties of each item.

In IRT we use an ability scale which may be thought of as representing the set of skills, abilities and knowledge that

contribute to performance. This scale is calibrated to have a mean of zero and ranges (theoretically) from negative to positive infinity. The units of ability are known as 'logits', where a logit is given by equation 4.

$$\text{logit}[P(\theta)] = \exp[ka(\theta - b_j)]$$

Equation 4. Definition of the logit – the unit of ability in psychometrics.

Item difficulty

For a dichotomous (two-category) item (yes or no; right or wrong, etc.), item difficulty is defined as the point on the measurement scale at which the probability of success is 0.5. For a polytomous item that carries several possible grades (usually the case for NCEA and tertiary examinations), we must estimate a difficulty parameter for each available grade, except the lowest.

Item discrimination

Item discrimination is the gradient of the item characteristic function at the point at which the probability of correct response is 0.5 (i.e. the value of the derivative of the function at this point), and theoretically can range between zero and infinity. The steeper the curve, the more highly the item discriminates between candidates of differing abilities, because, when the value of the gradient is high, small variations in ability give rise to significant differences in the probability of attaining a particular grade. However, very high discrimination values are undesirable for the same reason that very high item-total correlations are undesirable; they indicate redundancy amongst items. The ideal range for the discrimination is between about 1.0 and about 3.0. Table 3 shows the item parameters for the 2010 examination for Biology 90168.

Table 3. Difficulty and discrimination parameters for Biology 90168 under Samejima's Graded Response Model (1969).

Item	Discrimination	Difficulty (AME)	Difficulty (ME)	Difficulty (E)
Q1	1.51	-1.52	0.37	2.15
Q2	0.73	-2.57	0.59	4.68
Q3	2.43	-0.41	0.54	1.75

We see that all of the discrimination parameters of Table 3 fall within the desirable range. We also see that the items vary considerably in difficulty at each grade. In particular, it is relatively easy to obtain an *Achieved* grade or better in item 2 while for the same item it is very difficult to obtain *Excellence*.

Item characteristic curves

In IRT we depict graphically the performance of an item using item characteristic curves; plots showing the probability of achieving each available grade for an assessment as functions of candidates' ability. Figure 4 gives an example of a two-parameter item characteristic curve for an item that carries four grades, as is the case for NCEA external examinations and many examinations at tertiary level. The four curves represent the probabilities of achieving each grade for all candidates responding to the item. Each item in a given examination or test has its own unique set of characteristic curves.

The horizontal axis is the measurement scale on which candidates' abilities and item difficulties are estimated, and the vertical axis gives the probability of achieving a particular

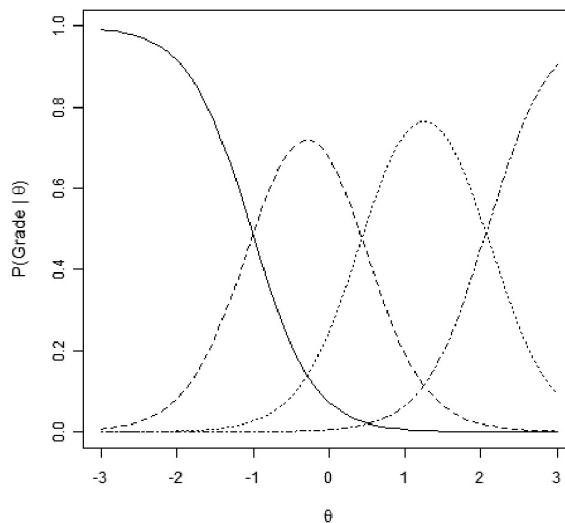


Figure 4. A typical set of item characteristic curves for NCEA external assessments constructed using Samejima's Graded Response Model (1969). The variable θ represents the measurement scale on which candidate ability and item difficulty are estimated.

grade. In this two-parameter item characteristic curve, and in equivalent plots later in this paper, the curve to the far left of the plot represents the probability of attaining a *Not Achieved* grade, and, moving left-to-right, the remaining curves represent the probabilities of attaining *Achieved*, *Merit*, and *Excellence*, respectively.

In implementing these models, we assume that we can characterise a candidate's performance with a single dimension. Of course, no examination actually measures just one cognitive construct, but often the skills or knowledge that we wish to measure are sufficiently strongly correlated that, statistically, they can be treated as representing a single dimension.

Figure 5 shows item response curves pertaining to the four items of the 2010 examination for the Level 2 Chemistry standard 90308 (Describe the nature of structure and bonding in different substances).

All four items discriminate well (as shown by the relatively steep slopes of the item characteristic curves), but items 1 and 3 discriminate the best of the four. For each item we see that there is a clearly defined domain of ability for which each grade is the most probable grade.

Grade thresholds

The threshold values for *Achieved*, *Merit*, and *Excellence* are defined as those locations on the ability axis at which results of *Achieved* and *Not Achieved*, *Merit* and *Achieved*, and *Excellence* and *Merit*, are, respectively, equally probable. Usually, we plot thresholds (values of θ_{NA} , θ_{AM} and θ_{ME}) on a dot chart, a particularly effective way of depicting grade thresholds. Figure 6 shows the threshold plot for the four items of the 2010 examination for the Level 2 Chemistry standard 90308 (Describe the nature of structure and bonding in different substances)

In this example none of the items are either particularly difficult or particularly easy. Additionally, the thresholds are reasonably (though not highly) consistent across the four items. There is no overlap between the domain in which the four *Achieved* thresholds fall, and that of the *Merit* grade. However,

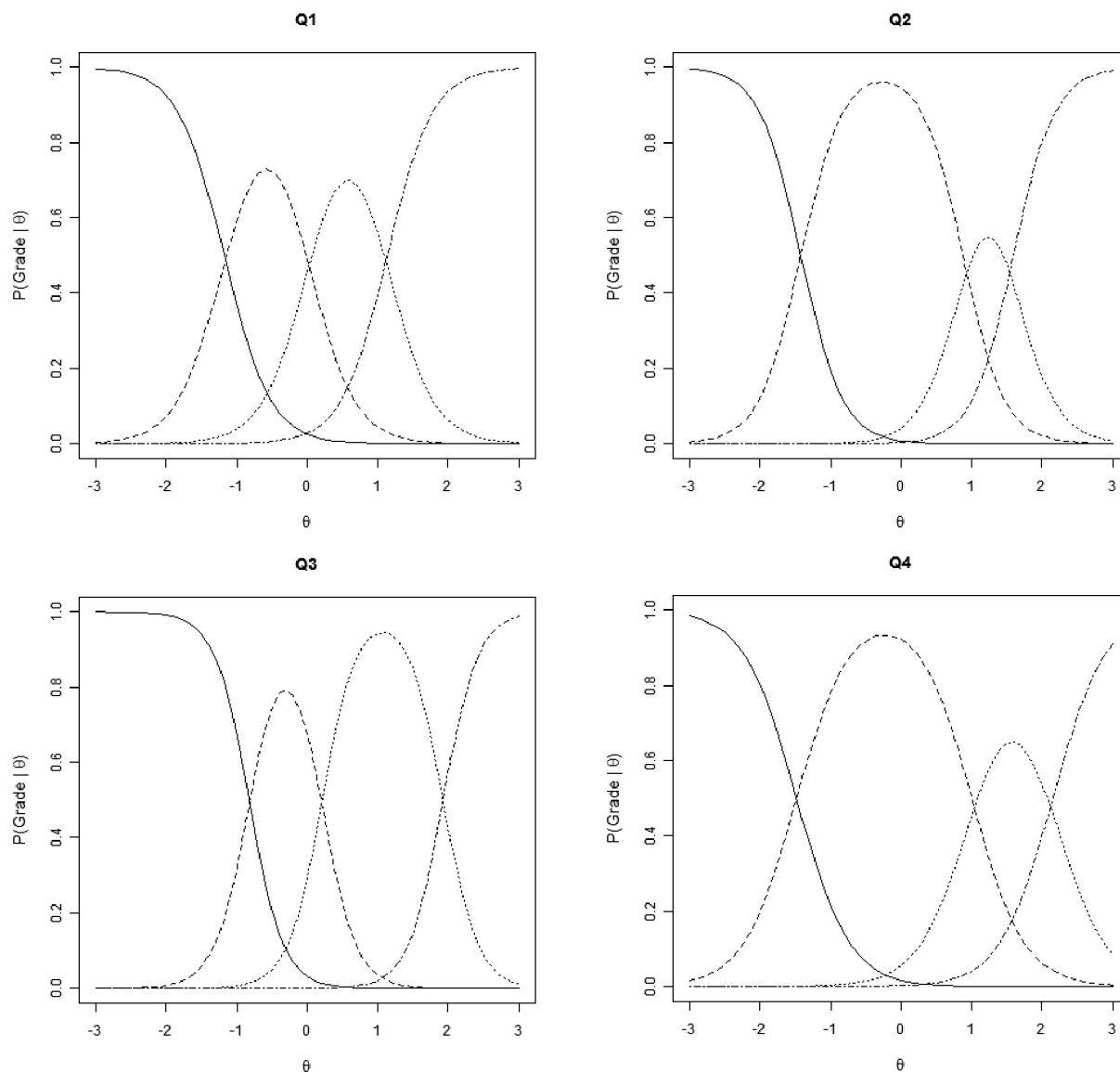


Figure 5. Item characteristic curves for the four items of the 2010 Level 2 Chemistry 90308 examination. From left to right the four curves represent the Not Achieved grade, the Achieved grade, the Merit grade and the Excellence grade. The variable θ represents the measurement scale on which candidate ability and item difficulty are estimated.

the *Merit* domain does overlap slightly with the *Excellence* domain; not a desirable property, although, in this case the overlap is not substantial.

Identifying item bias (differential item functioning)

Item bias, or differential item functioning (DIF), occurs when two or more groups of test or examination candidates, matched for overall ability, behave or perform differently on a particular item. We conduct DIF analysis in order to identify items that are possibly biased in favour of, or against, particular demographic groups (e.g. male or female candidates, or candidates identifying with different ethnic groups). Possibly, their different responses arise, not because one group of candidates has less knowledge of the subject matter, but because they held different assumptions initially or have had different cultural or other experiences.

During 2010 we developed analytic procedures for identifying DIF in NCEA assessments, based on those identified in the

literature (e.g. Zumbo 1999; 2007). We fit a series of ordinal logistic regression models to the results of groups of candidates that are matched for ability (e.g. males and females or students of different ethnicities). First, we fit a base ordinal logistic regression model (i.e. no covariates) to the set of item responses, then a regression with one covariate (e.g. group membership or gender). Finally, we fit more sophisticated models that include an interaction term (i.e. between ability and group membership or gender). These models are used to predict the item responses, where the main predictors are group membership and ability. For each model we calculate diagnostic statistics such as the log-likelihood and a Chi-square value (the log-likelihood for the base model minus the log-likelihood for each of the more complex models). Finally, the Chi-squared change for these models yields diagnostic statistics (i.e. the p-value and the R-squared change) which identify DIF. We detect the presence of DIF when the p-value is less than 0.05 and the R-squared change is greater than or equal to 0.035.

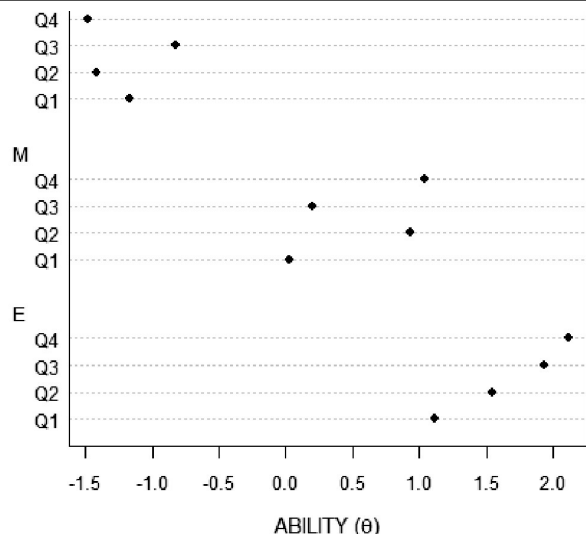


Figure 6. Grade thresholds for the four items (Q1 – Q4) of the 2010 Chemistry 90308 examination. Thresholds represent the points on the measurement scale at which adjacent categories are equally probable.

We may observe either uniform or non-uniform DIF. We have uniform DIF when one group has a higher probability of success on an item across the full range of abilities. We have non-uniform DIF when one group has a higher probability of success on an item on one or more domain of abilities, but has a lower probability on other domains. Our models produce output such as that of Table 4, pertaining to item 1 of the 2010 examination for the Geography standard 90704 (Select and apply skills and ideas in a geographic context).

The above item involved identifying particular geographic features on a satellite image and answering various questions that involved map reading skills. We see that this item exhibited uniform DIF between males and females (i.e. group membership was a significant predictor), but not between the ethnicity-based groups. Precisely why the item favoured males is not clear, but subject matter experts can often assist in such questions. It is important to note that the presence of DIF does not in itself establish bias. Bias is only established when the differential functioning is invalid in respect of the test construct, and the professional judgement of subject-matter experts is required to make this determination.

We can depict graphically the presence or otherwise of DIF. Figure 7 illustrates the presence of uniform DIF between male and female candidates for the above item. Note that the probability of success in the item is greater for male candidates than for female candidates across the entire ability domain.

To illustrate DIF we group the ability scores of all candidates in a set number of bins (here we use 12 bins, each of width 0.5 logits). We then plot

Figure 7. Graphical depiction of Differential Item Functioning for item 1 from the 2010 Geography standard 90704.

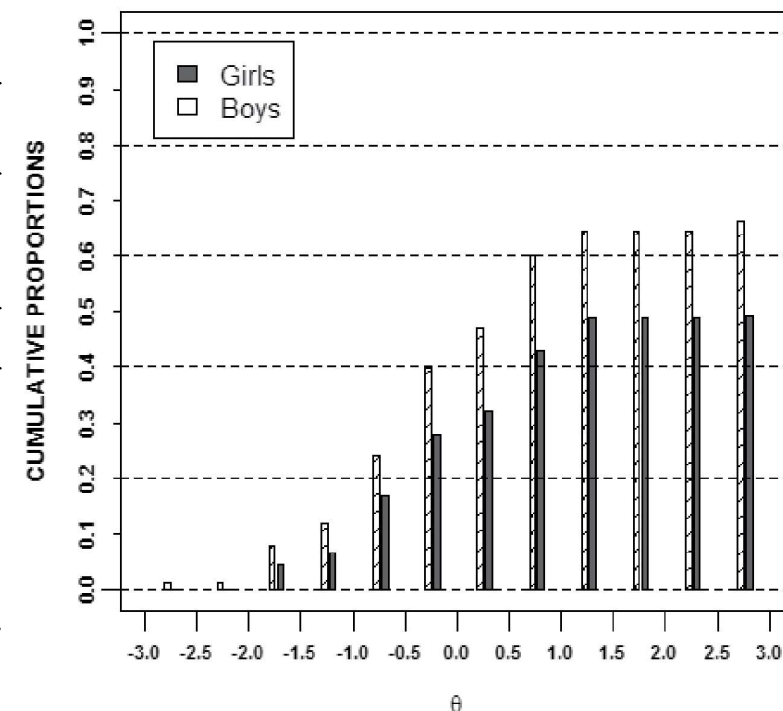


Table 4. Results of an analysis of Differential Item Functioning for item 1 from the 2010 examination for the Level 3 Geography standard 90704.

Comparison Groups	Uniform	Non-uniform
Male – Female	Yes	No
European – Māori	No	No
European – Pasifika	No	No
European – Asian	No	No

cumulative proportions of each subgroup (in this case males and females) attaining *Achieved*, *Merit* or *Excellence* grades (and whose estimated abilities fall within each bin), against the mean ability for each bin. For this particular item, across the entire domain of abilities, males were more successful than females. Nonetheless, our analyses of the results distributions of recent (i.e. the 2009 and 2010) examinations across many subjects and standards has revealed very little evidence of DIF.

New Zealand Scholarship: A hybrid of standards-based and normative assessment

New Zealand Scholarship examinations are designed to recognise high-level performance in a range of subjects (currently 35 subjects). Two passing grades are available for each subject: *Scholarship* and *Outstanding Scholarship*.

Results are awarded through a hybrid of normative assessment (in which candidates' grades depend on their performances relative to those of other candidates) and criterion-referenced assessment (in which candidates must satisfy established criteria for each available grade). In assessing candidates' scripts, each item is given a numerical (ordinal) score from 0 to 8, and the scores for individual items summed to produce an overall score for the script. Scores from 0 – 4 equate to a *No Award* grade; scores of 5 and 6 equate to a *Scholarship* grade, while scores of 7 and 8 equate to an *Outstanding Scholarship* grade.

Finally, a pair of cut scores, which define the range of total scores for award of Scholarship and Outstanding Scholarship for each script, is agreed. These cut scores are set so that about 3% of the NCEA Level 3 cohort, defined as the total number of candidates who have entered for 14 or more credits for NCEA Level 3 in that subject (not to be confused with the total number of students who have entered for the examination, which is usually a much smaller number), will receive a Scholarship, and about 0.4% will receive an Outstanding Scholarship. This is the normative part of the Scholarship assessment process.

Each script must include at least one item at Scholarship level if a Scholarship is to be awarded, and each script must include at least one item at Outstanding Scholarship level if an Outstanding Scholarship is to be awarded. If a script contains at least one item graded at 5 or 6, then we can say that the candidate has provided evidence of performance at Scholarship level, and similarly for Outstanding Scholarship level. This is the criterion-referenced part of the Scholarship assessment process.

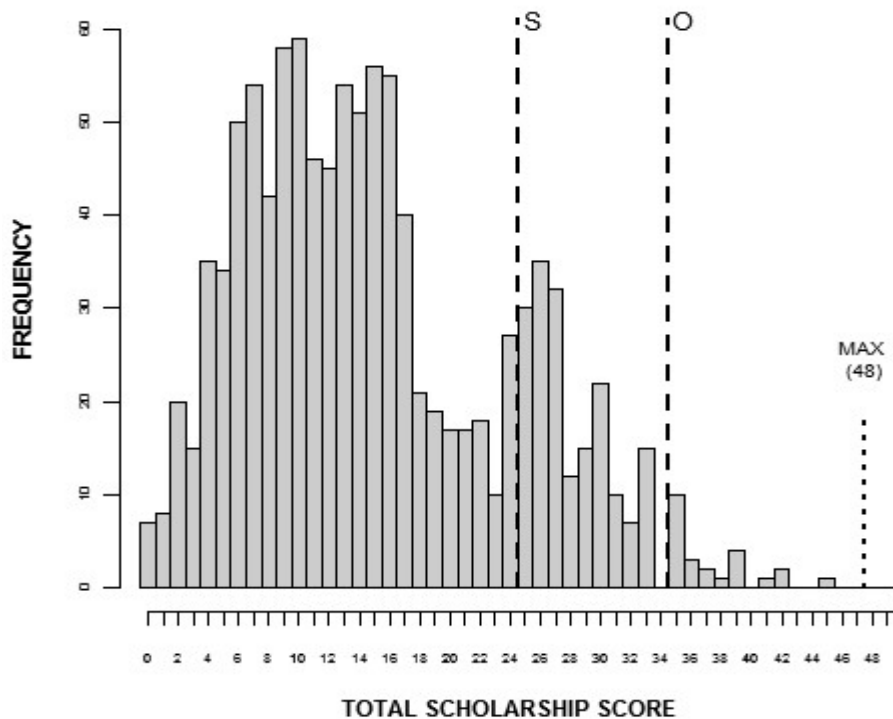
Awarding New Zealand Scholarship

Let us consider the 2010 Scholarship examination in Physics. This examination involved six items, so that the maximum possible score was 48. Following completion of the marking process, the cut score for Scholarship Physics was agreed at 25 (i.e. roughly 3% of the Physics Level 3 cohort) and the cut score for Outstanding Scholarship was set at 35 (roughly 0.4% of the cohort). Figure 8 gives a bar chart of total scores for the six-item 2010 Scholarship examination in Physics. The vertical lines indicate the cut scores for Scholarship (S) and Outstanding Scholarship (O) awards in that subject.

The bar chart shows a very wide range of performances on this examination. The Scholarship cut score of 25 was chosen so that roughly 3% of the cohort earned that score or above, all candidates at this score or above receiving at least one score of 5 over the complement of six items. The Outstanding Scholarship cut score of 35 was chosen so that roughly 0.4% of the cohort earned that score or above, all candidates at this score or above earning at least one score of 7.

The bar chart shows a highly skewed distribution of scores, a desirable characteristic in an examination that is designed to challenge top students. The positively-skewed distribution indicates that the test provides the most reliable information in the region of performance in which cut scores are likely to be set; around the midpoint of the total-score range for the *Scholarship* cut, and the three-quarters point for the *Outstanding Scholarship* cut.

Figure 8. Bar chart of total scores for the 2010 NZ Scholarship Physics examination, with Scholarship and Outstanding Scholarship cut scores (25 and 35 respectively).



Statistical modelling of New Zealand Scholarship

For all New Zealand Scholarship examinations we conduct similar analyses to those conducted for NCEA; dimensional analysis, IRT, etc, although the scholarship analyses are implemented on the full set of results, rather than on a sample. However, one additional analysis involves characterising the relationship between the results attained by Scholarship candidates in NCEA Level 3 in a given subject and their results in the Scholarship examination. Figure 9 gives a scatter-plot relating candidates' performances in the Level 3 Physics standards against their performances in Scholarship Physics. The vertical axis gives the mean expected percentiles (a measure of performance expressed as the expected percentile of the Level 3 candidature earned by the 'typical' candidate who has earned a particular grade in one of the external assessments) for each of the Level 3 Physics assessments taken by each candidate. The horizontal axis gives the total score earned by each candidate in the 2010 Scholarship Physics examination.

What exactly is a mean expected percentile? Let's illustrate using the Level 3 Physics examination for the four-credit Level 3 standard 90520 (Demonstrate understanding of wave systems). The national grade distribution for this examination was as follows: *Not Achieved* (24.0%), *Achieved* (54.5%), *Merit* (15.4%) and *Excellence* (6.2%). In the absence of precise information about any given student, our best estimate is that a student earning a *Not Achieved* grade sits at 12% of the candidature from the lowest score (i.e. the 12th percentile). Our best estimate is that a student earning an *Achieved* grade sits at 24% plus half of 54.5% (or the 51st percentile) from the lowest scoring candidate. Similarly, our best estimate is that a student earning a *Merit* grade sits at the 86th percentile, and a student earning an *Excellence* grade sits at the 97th percentile. Of course, each Scholarship candidate who took NCEA (some take other assessments such as Cambridge International Examinations or the International Baccalaureate) will have gained a particular set

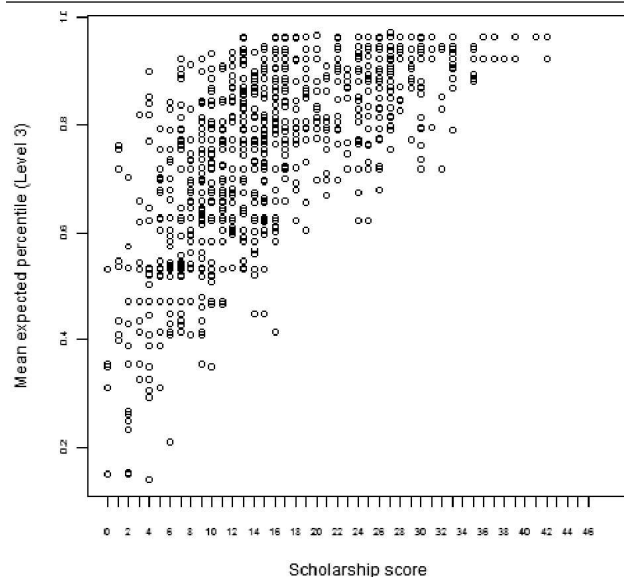


Figure 9. A plot of the mean expected percentiles for the 2010 Level 3 results for all Scholarship candidates in Physics against their total scores for the 2010 NZ Scholarship Physics examination.

of results in one or more of the four Level 3 Physics standards, and each is accorded a mean expected percentile for each of his or her Level 3 assessments. It is these percentiles, expressed as decimals, that are recorded on the vertical axis of Figure 9.

Essentially, this analysis illustrates the power of NCEA Level 3 in predicting performance in New Zealand Scholarship. In general terms the greater the mean expected percentile of the Level 3 assessments, the greater is the total Scholarship score. The relationship appears to be almost linear up to a Scholarship score of about 16, after which the curve levels off somewhat.

Figure 9 illustrates a particularly desirable attribute of a Scholarship examination: it extends the top end of performance of the Level 3 cohort. Students scoring in the top half of the Scholarship range typically achieve results at *Merit* and *Excellence* at Level 3. The examination has displayed discriminative power at higher levels of candidate performance than the Level 3 examinations.

Summary

Statistical modelling of NCEA and New Zealand Scholarship results provides very valuable feedback that supports ongoing improvement of our assessment processes. In addition to the

analyses described in this paper, we undertake many other diagnostic analyses that help to ensure fair and consistent assessment. Further applications of IRT are anticipated for the future. Eventually, our modelling programme will support the creation of banks of strongly-performing items for use by examiners and teachers, and in which we can have a very high degree of confidence.

It is important to be clear that the programme of analysis presented here is statistical in nature and concentrates on properties internal to the assessments themselves. The analyses we have described are necessary to ensure that assessments measure reliably, efficiently and fairly. They are not, however, of themselves sufficient to ensure valid measurement. Validity is the most essential property of any assessment and requires substantial content knowledge and understanding of the purposes of the assessment. Nonetheless, an assessment without strong reliability or that is of inappropriate difficulty, will not be valid, regardless of its specific content. Thus, the analyses described in this paper are essential for ensuring fair, reliable and valid national assessments for secondary-school qualifications in New Zealand.

References

- Cattell, R.B. 1966. The scree test for the number of factors. *Multivariate Behavioral Research* 1: 629–637.
- Cronbach, L.J. 1951. Coefficient alpha and the internal structure of tests. *Psychometrika* 16(3): 297–334.
- Hambleton, R.K., Swaminathan, H.; Rogers, H.J. 1991. *Fundamentals of Item Response Theory*. Newbury Park, California, Sage Publications. ISBN 0-8039-3647-8.
- Kaiser, H.F. 1960. The application of electronic computers to factor analysis. *Educational and Psychological Measurement* 20: 141–151.
- Pearson, K. 1901. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine* 2 (6): 559–572. <http://stat.smmu.edu.cn/history/pearson1901.pdf>.
- Samejima, F. 1969. Estimation of latent ability using a response pattern of graded scores. *Psychometric Monograph No. 17, Psychometric Society, Richmond, Virginia*.
- Zumbo, B.D. 1999. *A Handbook on the Theory and Methods of Differential Item Functioning (DIF): Logistic regression modelling as a unitary framework for binary and Likert-type (ordinal) item scores*. Ottawa, Directorate of Human Resources Research and Evaluation, Department of National Defense. Retrieved from <http://educ.ubc.ca/faculty/zumbo/DIF/handbook.pdf>
- Zumbo, B.D. 2007. Three generations of DIF analyses: Considering where it has been, where it is now, and where it is going. *Language Assessment Quarterly* 4(2): 223–233.

Acknowledgments

The authors wish to acknowledge the valuable contribution of Vernon Mogol to the work described in this paper, particularly in relation to Differential Item Functioning.