# Introduction to Classical Test Theory with CITAS

## ASSESSMENT SYSTEMS
### FOR GOOD MEASURE™

Nathan A. Thompson, Ph.D.
*Chief Product Officer*

# Table of Contents

# Contact Information

Assessment Systems Corporation
111 Cheshire Lane, #50
Minnetonka, MN 55305
*Voice:* 763.476.4764
*E-Mail:* solutions@assess.com
www.assess.com

Tests are used to make decisions on people every day.  In many cases, these are important decisions.  Unfortunately, most of the tests in the world are still of low quality.  They are not built according to international guidelines and best practices, and never evaluated for performance to check whether they are even producing reliable scores.

This white paper is intended for any individual that is interested in learning how to make tests and assessments better, by helping you apply international best practices for evaluating the performance of your assessments.  CITAS provides basic analytics necessary for this evaluation, and it does so without requiring advanced knowledge of psychometrics or of software programming.  However, if you are interested in more advanced capabilities and sophisticated psychometrics, I recommend that you check out www.assess.com/iteman for **Classical Test Theory** and www.assess.com/xcalibre for **Item Response Theory**.

What are the guidelines?  There are several resources, and can differ based on the use of the test as well as your country.  General guidelines are published by APA/AERA/NCME and the International Test Commission.  If you work with professional certifications, look at the National Commission for Certifying Agencies or the American National Standards Institute.  In the US, there are the Uniform Guidelines for personnel selection.

This paper will begin by defining the concepts and statistics used in classical item and test analysis, and then present how the CITAS spreadsheet provides the relevant information.  CITAS was designed to provide software for quantitative analysis of testing data that is as straightforward as possible – no command code, no data formatting, no complex interface.

# What is Classical Test Theory?

Why is it called "classical?"  This is to differentiate this type of analysis from the modern test analysis approach called *item response theory* (IRT) or sometimes *modern test theory*.  IRT is much more powerful, but only works with sample sizes numbering in the hundreds or larger.  This makes it extremely important in large-scale testing, but completely inappropriate for classroom-sized samples or other small-scale situations (N<100).

## *Classical analysis at the test level*

Classical test analysis is based primarily on number-correct (NC) or sum scores.  CITAS calculates descriptive statistics of the NC scores, as well as two important indices from classical test theory: KR-20 ($\alpha$) reliability, and the standard error of measurement (SEM).

*Reliability* is a classical concept that seeks to quantify the consistency or repeatability of measurement.  If a test is producing consistent scores, then we say it is reliable.  As to whether the scores actually mean what is intended, is part of a larger and more difficult issue called *validity*.  What is meant by consistency?  Let us assume that a student has a *true score* of 44 items out of 50.  If they took the test multiple times (assuming that we wiped their memory of

the test), they might get a 45, or a 43, etc.  This is consistent.  If they were to score a 34, then a 47, then a 39, it would be an unreliable test.

There are several approaches to indexing reliability, the most common of which is *internal consistency* using the Kuder-Richardson Formula 20 (KR-20) index.  Another index called *coefficient α* is also common, but is nothing more than a generalization of KR-20 to *polytomous* (rating scale or partial credit) data.

KR-20 ranges in theory from 0.0 to 1.0, with 0.0 indicating random number generation and 1.0 indicating perfectly consistent measurement.  In rare cases, it can even dip below 0.0.  Therefore, a higher number is regarded as better.  However, KR-20 is partially a function of test length, and tends to be higher when the test has more items.  So for shorter tests like 20 or 50 items, it is unrealistic to expect KR-20 values near 1.0.  In some cases, a value of 0.7 might be sufficient.

KR-20 is also important because it is used in the calculation of the SEM.  The SEM takes the concept of measurement consistency and applies it to student scores.  If we take plus or minus 2 SEMs around a student's observed score, that gives us a range we are 95% confident contains their true score.  If this number is very small, this means that we have an accurate estimate of the true score.  If it is large, we do not have an accurate estimate.  Suppose the SEM is 5.0.  Then the interval is plus or minus 10.  For a student with a test score of 39, this means we expect their true score to be anywhere from 29 to 49:  hardly an accurate test!

## *Classical analysis at the item level*

The statistics presented so far only provide information at the level of the entire test.  Classical test theory also has statistics for evaluating individual items from a quantitative perspective.  The goal of item analysis is to use detailed statistics to determine possible flaws in the item, and then decide whether to revise, replace, or retire the item.  This can be something as specific as identifying a bad distractor because it pulled a few high-ability examinees, or something as general as "this item is harder than I prefer for my students."

In large-scale or high-stakes testing, item analysis is typically performed before the test goes "live" to ensure only that quality items are used.  Often, it is done after pretesting the items on some small set of the population, perhaps by inserting them as unscored into the test the year before they are to be used as scored items.  In cases where pretesting is not feasible, item analysis (as well as equating and standard-setting) can take place after the live administration.  However, this requires that scores be reported much later, in some cases weeks or months later.  This is of course unacceptable for classroom assessment, but CITAS allows you to evaluate items immediately after test delivery.

Item analysis is important because it is analogous to quality control of parts used in the assembly of a final manufacturing product.  Nobody wants bad tires or bad brake pads in their car.  Releasing tests with bad or untested items is like releasing cars off the assembly line with bad or untested brake pads.  Both prevent plenty of opportunities for litigation.

In classical test theory, there are primarily two concepts we are interested in assessing: item difficulty and item discrimination.  *Item difficulty* is a simple concept in classical test theory; it simply refers to the proportion of students that correctly answered an item.  This is called the *P*-value.  Yes, I know, we already use "*p*-value" as the term for statistical significance, but it's too

late to change a 100-year-old theory.  If 95% correctly answered it, the item is quite easy.  If 30% correctly answer, then the item is quite difficult, especially when you consider that a four-option multiple choice item presents a 25% chance of guessing the correct answer!  For this reason, items with $P < 0.50$ are generally considered to be quite difficult, while we typically see them more in the 0.70 to 0.80 range.  However, specific tests might produce different ranges of statistics, requiring you to shift the paradigm somewhat.

Just what is too difficult or too easy?  That is a judgment call that you have to make while taking into account the content of the item, the purpose of the test, and the sample of students.  A test that is designed to be extremely difficult might have an average $P$ of 0.60.  Conversely, a test that is given to a group of extremely high ability examinees can be expected to have an average $P$ of 0.90.  Regardless of the average $P$, it is often preferable to have items with a range of difficulty.  If you have no items with $P < 0.70$, it means that all the items were fairly easy, and there was not a single item on the test to "separate the men from the boys."  This might be acceptable if the purpose of the test is just to assess entry level knowledge to a topic, but if the purpose of the test is to identify the top students, a test with all item $P$ values above 0.70 would not do a good job.

Because the scale is conceptually reversed (an easy item has a high number, and a difficulty item has a low number), some textbooks refer to this as *item facility*.  But this is a minority.

*Item discrimination* refers to the power of the item to differentiate between examinees with high and low levels of knowledge or ability.  But of course we never know the true score of any examinee.  The best estimate we have is the total NC score on the test.  Therefore, item discrimination is typically defined as the correlation between item scores (scored 0/1) and total test scores, called the *item-total correlation*.  The equation used to calculate this is called the *point-biserial correlation*, or $r_{pbis}$, though some researchers prefer to use its cousin called the *biserial correlation*.  This provides an index of whether students who get the item correct are scoring highly, which is the hallmark of a good item.

The item-total correlation has the advantage that it is interpreted as any other correlation, which many people are used to working with.  An $r_{pbis}$ of 0.0 indicates that there is no correlation, which means that there is no relationship between the item and the total scores.  This means that the item is providing no information, and item responses are essentially random with respect to total scores.  But as $r_{pbis}$ increases, it indicates a stronger relationship between the item and total score.  A value of 0.20 means a decent item, and highly discriminating items will have values in the 0.50 or 0.60 range.

On the other hand, a negative $r_{pbis}$ is very bad news.  This means that there is an inverse relationship, namely that low-scoring students performed better on the item than high-scoring students did.  This typically indicates one of three things:

1. A key error;
2. A very attractive distractor;
3. This item is so easy/hard that there are few examinees on one side of the fence, making it difficult to correlate anything.

All three things are issues with the item that need to be addressed.

For completeness' sake I must mention another statistic, sometimes called the *classical difficulty index* or the *top-bottom index.* This was developed before the $r_{pbis}$ but is occasionally still used. It is based on the same concept that we want high-ability students to get the item correct more often than low-ability students. So we divide the sample in half, and find the proportion of the top half correctly answering and subtract the proportion of the bottom half correctly answering (sometimes done with the top and bottom 27%). Like the $r_{pbis}$, a positive value indicates a better quality item.

A very important thing to note about the item statistics: like all statistics, their stability depends on sample size. In general, we need 20 or 30 people to get marginally useful statistics, and they start becoming statistically stable near 50 examinees. Therefore, while CITAS results with 20 students will provide some helpful information, do not consider the item statistics to be perfectly stable.
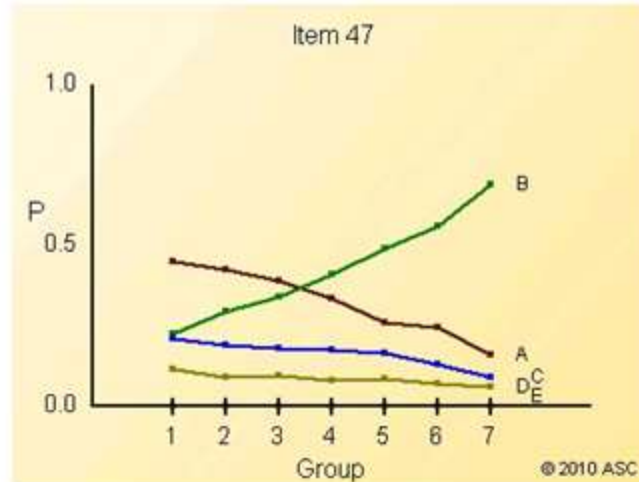
## *Option Statistics*

If you wish to dig even deeper into the performance of an item, the next step is an evaluation of *option statistics*. With multiple choice items, the word *option* refers to the possible answers available. The correct answer is typically called a *key* and the incorrect options called *distractors*.

Evaluating the option statistics for telltale patterns is an important process in diagnosing items that have been flagged for poor *P* or $r_{pbis}$ values at the item level. This is done by evaluating *P* and $r_{pbis}$ at the option level. In general, we want two things to happen:

1. The P for the key is greater than the P for any of the distractors. That is, we don't want more students choosing one of the distractors than the key. In many (but not all) cases, this means the distractor is arguably correct or the key is arguably incorrect.

2. The $r_{pbis}$ for the distractors should be negative but the $r_{pbis}$ for the key should be positive. If an $r_{pbis}$ for a distractor is positive, this means that smart examinees are choosing it, and we usually want the not-so-smart examinees selecting the incorrect answers. However, this pattern is very susceptible to fluctuations in small sample sizes; if only 4 examinees select an option and one or two are of very high ability, that is often enough to produce a positive $r_{pbis}$ and therefore flag the item.

An even deeper analysis is called *quantile plots*. This methodology splits the sample into quantiles based on ability (total score) and then evaluates the option P values for each, plotting them on a single graph. We want to look for the same pattern of performance, which usually means that the line for the key has positive slope (positive $r_{pbis}$) and that the lines for the distractors have negative slope (negative $r_{pbis}$). An example of this, from *Iteman*, is below.

Item 47

**Item information**

| Seq. | ID | Key | Scored | Num Options | Domain |
|------|----|----|--------|-------------|--------|
| 47 | 47 | B | Yes | 5 | 1 |

**Item statistics**

| N | P | Total Rpbis | Total Rbis | Alpha w/o | Flags |
|---|---|-------------|------------|-----------|-------|
| 16658 | 0.447 | 0.307 | 0.387 | 0.804 | |

**Option statistics**

| Option | N | Prop. | Rpbis | Rbis | Mean | SD | Color | |
|--------|---|-------|-------|------|------|----|-------|--|
| A | 5199 | 0.312 | -0.214 | -0.280 | 27.035 | 6.396 | Maroon | |
| B | 7448 | 0.447 | 0.307 | 0.387 | 31.455 | 6.237 | Green | **KEY** |
| C | 2621 | 0.157 | -0.102 | -0.154 | 27.586 | 6.445 | Blue | |
| D | 1373 | 0.082 | -0.057 | -0.103 | 27.896 | 6.826 | Olive | |
| E | 3 | 0.000 | -0.010 | -0.189 | 24.333 | 8.083 | Gray | |
| Omit | 14 | 0.001 | -0.011 | -0.113 | 25.214 | 7.748 | | |
| Not Admin | 14 | | | | 25.214 | 7.748 | | |

**Quantile plot data**

| Option | N | 0-14% | 14-29% | 29-43% | 43-57% | 57-71% | 71-86% | 86-100% | Color | |
|--------|---|-------|--------|--------|--------|--------|--------|---------|-------|--|
| A | 5199 | 0.448 | 0.426 | 0.390 | 0.335 | 0.260 | 0.244 | 0.158 | Maroon | |
| B | 7448 | 0.227 | 0.294 | 0.338 | 0.411 | 0.489 | 0.558 | 0.692 | Green | **KEY** |
| C | 2621 | 0.210 | 0.188 | 0.178 | 0.175 | 0.167 | 0.129 | 0.089 | Blue | |
| D | 1373 | 0.113 | 0.091 | 0.095 | 0.080 | 0.085 | 0.068 | 0.061 | Olive | |
| E | 3 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | Gray | |

## What is CITAS?

The Classical Item and Test Analysis Spreadsheet (CITAS) is a simple tool to statistically analyze small-scale assessments, available at www.assess.com/citas.  CITAS is a Microsoft Excel® spreadsheet with all necessary calculations programmed in as formulas, which means that all the user is required to do it type or paste in the student responses and the correct answers, or keys. CITAS will then score all students with number-correct (NC) scoring, as well as populate

important statistics.  Statistics include both test-level statistics such as reliability, and item-level statistics such as difficulty ($P$) and discrimination ($r_{pbis}$).  CITAS statistics are listed in Table 1.

*Table 1: Statistics in CITAS output*

| Test-level statistics | Item statistics |
|---|---|
| Number of examinees | $P$ |
| Number of items | $r_{pbis}$ |
| NC score mean | Number correct |
| NC score standard deviation | Number incorrect |
| NC score variance | Mean score correct |
| Minimum NC score | Mean score incorrect |
| Maximum NC score | Response frequencies |
| KR-20 (alpha) reliability | Response mean scores |
| Standard error of measurement (SEM) | |
| Mean $P$ | |
| Minimum $P$ | |
| Maximum $P$ | |
| Mean $r_{pbis}$ | |
| Minimum $r_{pbis}$ | |
| Maximum $r_{pbis}$ | |

## Classical test analysis with CITAS

CITAS provides both the KR-20 and SEM, as well as simple descriptive statistics of the student scores.  Table 1 explains the summary statistics found on the "Output" tab

*Table 2: Summary test-level statistics in CITAS output*

| Summary Statistic | Definition |
|---|---|
| *Test-level* | |
| Examinees: | Number of students |
| Items: | Number of items |
| Mean: | Average NC test score |
| SD: | Standard deviation of NC scores |
| Variance: | Variance of NC scores |
| Min: | Lowest score |
| Max: | Highest score |
| KR-20: | Reliability of measurement |
| SEM: | Standard error of measurement |

Example output is shown in Table 3. This test has 72 items, which were answered by 100 students. The test was quite fairly difficult, with an average score of 40.32 out of 72. There was a wide spread of scores, with an SD = 11.01 and a range of 14 to 69. The test had sufficient reliability, with a KR-20 of 0.89 and SEM of 3.67.

*Table 3: Summary test-level statistics in CITAS output*

| Statistic | Value |
|---|---|
| Examinees: | 100 |
| Items: | 72 |
| Mean: | 40.32 |
| SD: | 11.01 |
| Variance: | 121.29 |
| Min: | 14 |
| Max: | 69 |
| KR-20: | 0.89 |
| SEM: | 3.67 |

## Classical item analysis with CITAS

CITAS presents the $P$ and $r_{pbis}$ for each item as well as some supplemental statistics based on the correct/incorrect dichotomy. The first is the number of correct and incorrect responses. This is obviously a repackaging of the $P$ value, but provides an alternative method of looking at difficulty if you prefer to use it. Additionally, CITAS presents the mean scores for students who got the item correct and incorrect. If the item is discriminating well, the mean score will be higher for the "correct" students. Similarly, this is a repackaging of the $r_{pbis}$, but provides an alternative method of evaluating item discrimination.

Table 4 presents example results from a test of 20 items. Let us go through the results for the first five items and interpret the statistics.

Item 1 is a fairly easy item with a $P$ of 0.942, and has a minimal but still positive $r_{pbis}$ at 0.18. This positive discrimination is reflected in the mean scores; the average score for examinees responding correctly is higher than examinees responding incorrectly.

Item 2 is similar to Item 1, but is more discriminating, with an $r_{pbis}$ of 0.33 and wider point difference. This item is a solid item for the test even though it is relatively easy.

Item 3 presents an ideal item from the classical perspective. It is more difficult than the first two items, but still not all that difficult in an absolute sense with 61% of the examinees responding correctly. More importantly, it has a very strong discrimination; the $r_{pbis}$ is 0.49 and a wide point difference. Items like this are very powerful, as we can achieve decent test reliability with only a few items.

Item 7 is an item that should be flagged and reviewed. While 77% answered correctly, the discrimination is negative, indicating that an option was probably too strong of a distractor.

Item 8 represents a very difficult item, yet strong. Only 40% of students responded correctly, but the discrimination is still 0.49.

*Table 4: CITAS item statistics*

| Item | P | Rpbis | Number correct | Number incorrect | Mean score correct | Mean score incorrect |
|------|------|-------|----------------|------------------|--------------------|----------------------|
| 1 | 0.94 | 0.18 | 94 | 6 | 40.84 | 32.50 |
| 2 | 0.89 | 0.33 | 89 | 11 | 41.63 | 29.91 |
| 3 | 0.61 | 0.49 | 60 | 39 | 44.78 | 33.49 |
| 4 | 0.81 | 0.37 | 81 | 19 | 42.33 | 31.84 |
| 5 | 0.82 | 0.34 | 82 | 18 | 42.10 | 32.33 |
| 6 | 0.58 | 0.43 | 58 | 42 | 44.38 | 34.76 |
| 7 | 0.77 | -0.08 | 76 | 23 | 39.83 | 41.91 |
| 8 | 0.40 | 0.49 | 40 | 60 | 47.08 | 35.85 |
| 9 | 0.62 | 0.45 | 61 | 38 | 44.31 | 33.95 |
| 10 | 0.90 | 0.22 | 90 | 10 | 41.17 | 32.90 |

## Distractor analysis with CITAS

In addition to evaluating items as a whole, statistics can be used to evaluate individual *options* of items.  The statistics for the correct option serve as the statistics for the item as a whole, because 90% of the students answered the correct answer of "A" then 0.90 is both the proportion of students who answered "A" and the proportion of students who answered correctly.  But what makes option statistics useful is the evaluation of the incorrect options, known as *distractors*.  This provides even greater detail about the performance of the item, as we will see.

CITAS provides a method to evaluate distractors by presenting the average scores for examinees with a given response.  If the item is performing well, examinees that respond correctly will have the highest average score.  Examinees responding incorrectly will have a lower average score.  When examining individual options, the option that is the most incorrect should have the lowest average score.  For example, if "A" is correct, "B" and "C" are incorrect, and "D" is not even close, then we would expect student who selected "A" to have high scores, and students who selected "D" to have low scores.

The final tab of CITAS presents statistics for distractor analysis, as seen in Table 5.  We saw in Table 4 that 94 examinees responded correctly to the first item, while 6 responded incorrectly.  Table 5 shows us that of those 6 incorrect responses, 3 chose "A," 1 chose "B," and 2 chose "D."

*Table 5: CITAS distractor analysis for example form 1*

| Item | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Key | C | B | B | B | D | C | B | A | C | B |
| **Option N** | | | | | | | | | | |
| A | 3 | 3 | 16 | 7 | 12 | 7 | 2 | 40 | 14 | 3 |
| B | 1 | 89 | 60 | 81 | 0 | 8 | 76 | 5 | 10 | 90 |
| C | 94 | 5 | 9 | 6 | 6 | 58 | 13 | 18 | 61 | 2 |
| D | 2 | 3 | 14 | 6 | 82 | 27 | 8 | 37 | 14 | 5 |
| E | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| | | | | | | | | | | |
| **Option P** | | | | | | | | | | |
| A | 0.03 | 0.03 | 0.16 | 0.07 | 0.12 | 0.07 | 0.02 | 0.40 | 0.14 | 0.03 |
| B | 0.01 | 0.89 | 0.60 | 0.81 | 0.00 | 0.08 | 0.76 | 0.05 | 0.10 | 0.90 |
| C | 0.94 | 0.05 | 0.09 | 0.06 | 0.06 | 0.58 | 0.13 | 0.18 | 0.61 | 0.02 |
| D | 0.02 | 0.03 | 0.14 | 0.06 | 0.82 | 0.27 | 0.08 | 0.37 | 0.14 | 0.05 |
| E | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| | | | | | | | | | | |
| **Option Rpbis** | | | | | | | | | | |
| A | 0.02 | -0.13 | -0.33 | -0.30 | -0.22 | -0.18 | 0.15 | 0.49 | -0.26 | -0.05 |
| B | -0.20 | 0.33 | 0.49 | 0.37 | ##### | -0.19 | -0.08 | -0.26 | -0.24 | 0.23 |
| C | 0.18 | -0.24 | -0.21 | -0.05 | -0.25 | 0.43 | -0.02 | -0.19 | 0.45 | -0.06 |
| D | -0.19 | -0.17 | -0.17 | -0.23 | 0.34 | -0.25 | 0.08 | -0.23 | -0.16 | -0.24 |
| E | ##### | ##### | ##### | ##### | ##### | ##### | ##### | ##### | ##### | ##### |

Here, we can see the general pattern of the key having a strong (positive) $r_{pbis}$ while the distractors have negative $r_{pbis}$. The exception, of course, is Item 7 as discussed earlier. There, the key (B) has $r_{pbis}$ =-0.08 while both A and D have positive $r_{pbis}$, albeit with small N. This means that those two distractors happened to pull some smart students, and should be reviewed.

In some cases, you might see a straying from the desired pattern. In Item 1, we see that A has a positive but small $r_{pbis}$ (0.02). There are only 3 examinees that selected A, so this is a case of the aforementioned situation where it takes only one smart examinee to select a distractor and it might be flagged as a positive $r_{pbis}$. This item is likely just fine.

## Summary statistics of item statistics

CITAS provides one additional level of information: summary statistics of the individual item statistics. These are presented in the lower right of the "Output" tab. Example results are show in Table 6. We again see that this is an fairly difficult test; the average *P* value was only 0.56, meaning that the average score was 56% correct. The most difficult item had a *P* of 0.18 and the easiest item had a *P* of 0.94.

The discrimination power of this test was quite good. The average $r_{pbis}$ was 0.34, and the lowest was item 5 with a value of -0.09.

*Table 6: Summary statistics in CITAS output*

| Summary Statistic | Value |
| --- | --- |
| Mean P: | 0.56 |
| Min P: | 0.18 |
| Max P: | 0.94 |
| Mean Rpbis: | 0.34 |
| Min Rpbis: | -0.09 |
| Max Rpbis: | 0.59 |

## Summary

Item analysis is a vital step in the test development cycle, as all tests are composed of items and good items are necessary for a good test.  Classical test theory provides some methods for evaluating items based on simple statistics like proportions, correlations, and averages.  However, this does not mean item evaluation is easy.  I've presented some guidelines and examples, but it really comes down to going through the statistical output and a copy of the test with an eye for detail.  While psychometricians and software can always give you the output with some explanation, it is only the item writer, instructor, or other content expert that can adequately evaluate the items because it requires a deep understanding of test content.

Although CITAS is quite efficient for classical analysis of small-scale assessments and teaching of classical psychometric methods, it is not designed large-scale use.  That role is filled by two other programs, *FastTest* and *Iteman 4*.  *Iteman 4* is designed to produce a comprehensive classical analysis, but in the form of a formal MS Word report ready for immediate delivery to content experts; please visit www.assess.com/iteman to learn more.  FastTest is ASC's comprehensive ecosystem for test development, delivery, and analytics.  It can produce Iteman reports directly from the system if you utilize it to deliver your tests.

## Further reading

Downing, S.M., & Haladyna, T.M. (Eds.) (2006). *Handbook of test development*.  Philadelphia: Taylor & Francis.

Furr, R.M., & Bacharach, V.R. (2007). *Psychometrics: An introduction*. Thousand Oaks, CA: Sage.

Shultz, K.S., & Whiney, D.J. (2005). *Measurement theory in action*. Thousand Oaks, CA: Sage.