# Reporting and Analytics - Is My Test Working?

Jim Parry, M.Ed., CPT, Compass Consultants, LLC
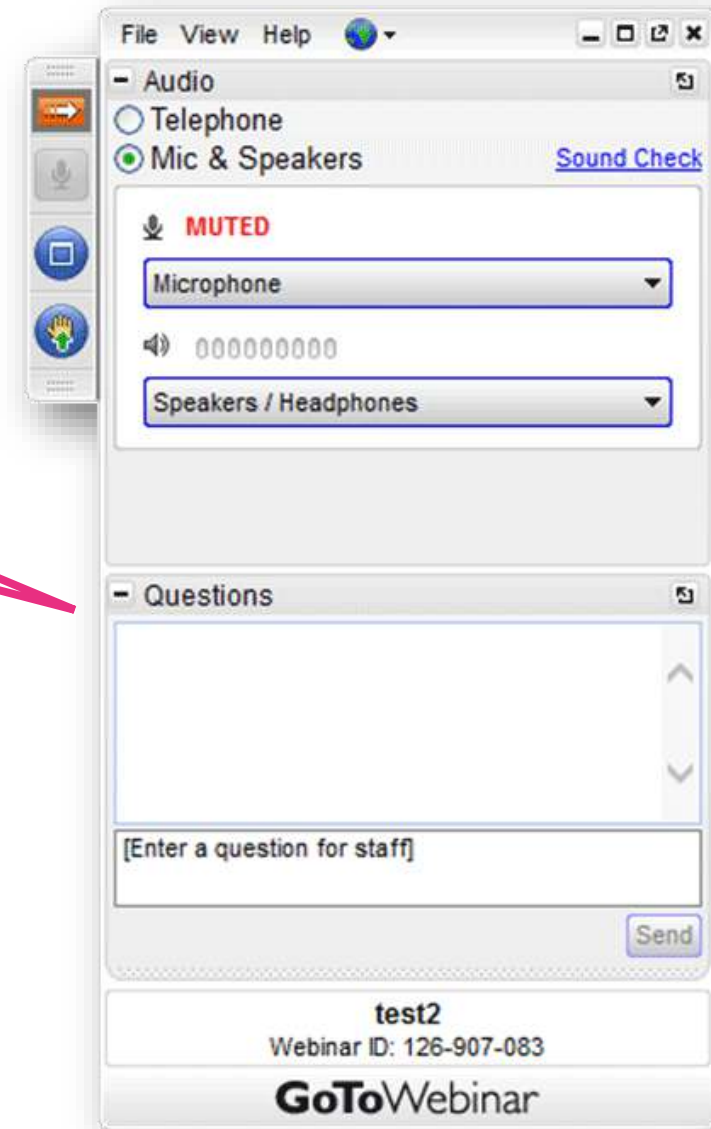
**May 26, 2021**

To ask questions,
use the "Questions"
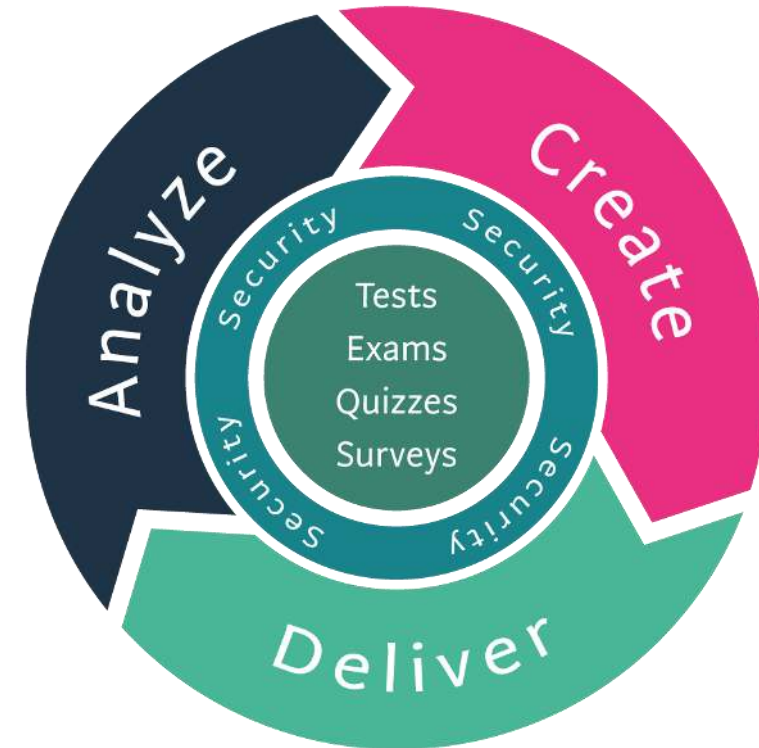feature

**Watch for an email after the webinar:**
- Download slides (PDF)
- View a recording
- Answer a survey

# About Questionmark

## Background

- Founded in 1988

- Assessment solutions to measure knowledge, skills, abilities and attitudes securely for certification, regulatory compliance, workforce learning, sales-force readiness and higher education

- ISO/IEC 27001 Certified (Learn more: www.questionmark.com/trust)



- *Questionmark OnDemand*
- *Questionmark OnDemand for Government*
- *Questionmark OnPremise*

# Today's Presenter

## Jim Parry, M.Ed., CPT, Compass Consultants, LLC

- Owner and Chief Executive Manager of Compass Consultants, LLC

- Over 40 years' experience in course design, development, presentation and assessment design and analysis

- Holds a Master of Education degree from the University of West Florida and is a Certified Performance Technologist (CPT), awarded by the International Society of Performance Improvement (ISPI)

- Has been presenter of pre-conference workshops and educational sessions at various professional conferences for many years

- Internationally recognized consultant providing services concerning test design, development, establishment of cut scores, and analysis

- Jim is a consulting partner of Questionmark

# About Compass Consultants, LLC

## Background

- Founded in 2010
- A leader in the application of Human Performance Technology (HPT), specializing in the design, development and presentation of training interventions and the psychometrics of test development and analysis.
- Learn more: www.gocompassconsultants.com

# Agenda

The Purpose of Test and Test-item Analysis

Commonly Reported Statistics

Item and Assessment Analysis

Reporting and Interpretation

# Legal Disclaimer

- The presentation may include information about legal issues and legal developments.  Such materials are for informational and/or educational purposes only and may not reflect the most current legal developments.  These informational/educational materials are not intended, and should not be taken, as legal advice on any particular set of facts or circumstances.  You should contact an attorney for advice on specific legal problems or questions.

- Information and/or software tools are provided "as is" without any express or implied warranty of any kind including warranties of merchantability, noninfringement of intellectual property, or fitness for any particular purpose. In no event shall Compass Consultants, LLC., or its employees, contractors, sub-contractors, agents, officers or attorneys be liable for any damages whatsoever (including, without limitation, damages for loss of profits, business interruption, loss of information) arising out of the use of or inability to use the information, even if Compass Consultants, LLC has been advised of the possibility of such damage.

# The Purpose of Test Analysis

Why Should I Care?

# question mark

# Quick Poll ☑

**How does your organization analyze tests and test-items?**

A. We are primarily concerned with the average or mean score.

B. We are primarily concerned with how many pass and how many fail.

C. We perform complete statistical analysis using Item Response Theory (IRT) or Classical Test Theory (CTT)

D. We do not perform any type of test or test-item analysis

# The Purpose of Test Analysis

- Provides feedback from test
  - Assists in improving test
- Ensures test items are valid and reliable
- Provides accurate measure of learner's output
- Can be used to pinpoint weak instruction
- Key piece of defensibility

# Classical Test Theory (CTT) vs. Item Response Theory (IRT)

| Classical Test Theory | Item Response Theory |
|---|---|
| • Longer tests are more reliable than shorter tests | • Shorter test are more reliable than longer tests |
| • Focuses on overall test performance | • Focuses on item performance |
| • Evaluates consistency among administrations of the test | • Takes into account the difficulty of items when estimating a test-takers ability |
| • Anything consistent is considered as the 'truth' even if test-taker consistently 'cheats' | • Performance of each item calculated individually after each administration |
| • Person ability depends on the test | • In theory, test-takers and item parameters are independent of each other – a person should have the same ability no matter which set of test-items they take |
| • Item parameters (difficulty and discrimination) depend on test takers | • An item should have the same difficulty and discrimination no matter who is taking the test |
| • About 100 examinees may be required to obtain stable results | • 500 – 1000 examinees may be required to obtain stable results |

# Commonly Reported Test Statistics

Numbers, Numbers, Numbers...
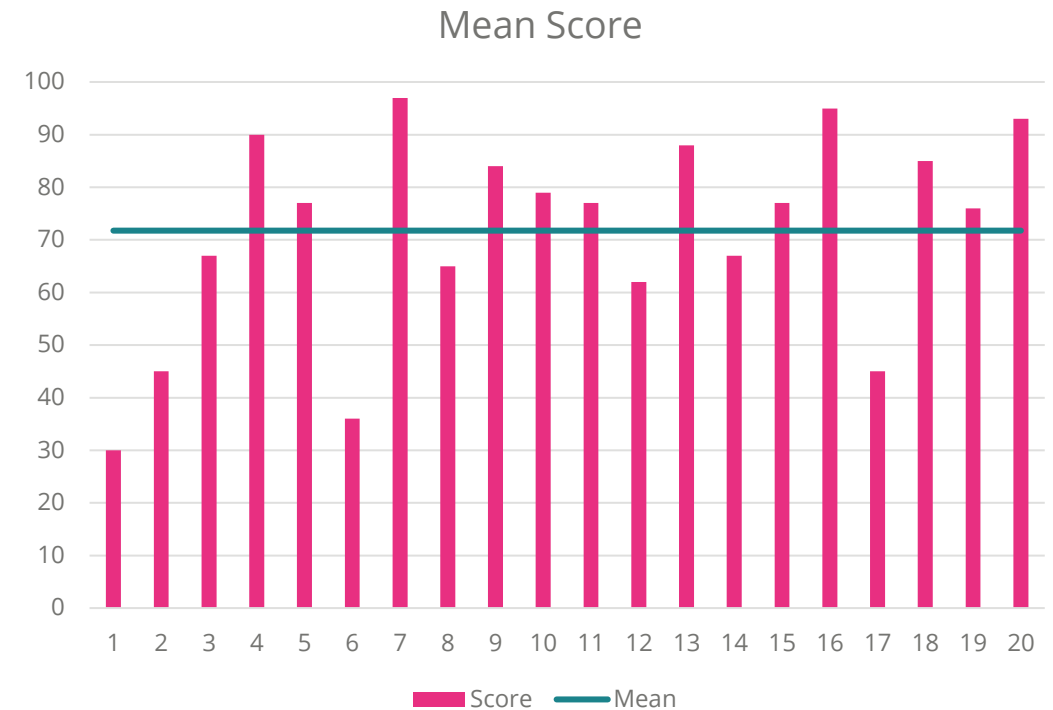
# Test Analysis Report (Questionmark Classic Reports)

## Table of Test Statistics

| | | | | | |
|---|---|---|---|---|---|
| Number of examinees | 27 | Mean | 38.96/64 (60.88%) | Standard error of mean | 1.77/64 (2.77%) |
| Number of items | 40 | Median | 41/64 (64.06%) | Standard error of measurement | 4.13/64 (6.45%) |
| Maximum possible score | 64 | Mode | 45/64 (70.31%) | Skew | -3.038 |
| Minimum achieved score | 0/64 (0%) | Standard deviation | 9.21/64 (14.39%) | Kurtosis | 12.444 |
| Maximum achieved score | 51/64 (79.69%) | Variance | 84.88/64 (132.62%) | Test reliability (Cronbach's Alpha) | 0.799 |

*Reliability is most meaningful if all items cover the same subject area.*

# Mean

- The MEAN of a set of test results is the AVERAGE score

  o Raw number or percentage

    - 30, 45, 67, 90, 77, 36, 97, 65, 84, 79, 77, 62, 88, 67, 77, 95, 45, 85, 76, 93

      - *Total = 1435*

      - *20 scores Δ 1435 / 20 = 71.75 average or mean*
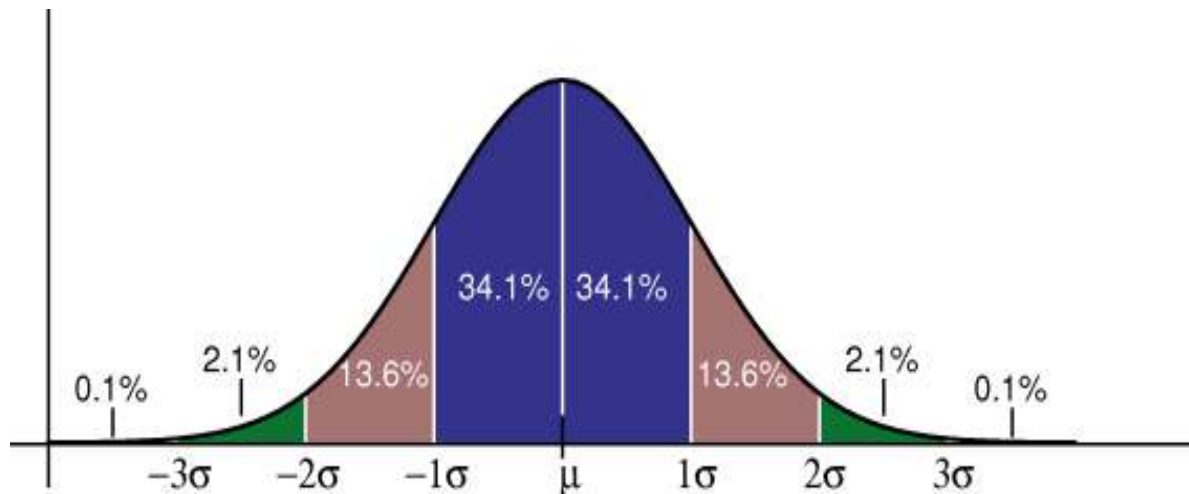
### Mean Score

# Median and Mode

- Median
  - The value lying at the midpoint of a distribution of numbers such that there is an equal probability of falling above or below it.
    - 30, 45, 67, 90, 77, 36, 97, 65, 84, 79, 77, 62, 88, 67, 77, 95, 45, 85, 76, 93
    - 30, 36, 45, 45, 62, 65, 67, 67, 76, **77, 77**, 77, 79, 84, 85, 88, 90, 93, 95, 97
      - *Median = 77*
  - May give a better indication of 'average' test performance if many scores were high or low

- Mode
  - The number which appears most often in a set of numbers – the 'peak' of a distribution where most samples concentrate
    - 30, 45, 67, 90, **77**, 36, 97, 65, 84, 79, **77**, 62, 88, 67, **77**, 95, 45, 85, 76, 93
      - *Mode = 77*

# Standard Deviation σ (SD)

- The standard deviation is a function of the bell curve that defines the average deviation or degree of distribution of scores from the mean score.

  o How far from the mean a sample of test takers deviate
  - 30, 45, 67, 90, 77, 36, 97, 65, 84, 79, 77, 62, 88, 67, 77, 95, 45, 85, 76, 93
    - *SD = 19.16*

| Standard Deviation | $\sigma$ = 19.162137 |
| --- | --- |
| Variance | $\sigma^2$ = 367.1875 |
| Count | n = 20 |
| Mean | $\mu$ = 71.75 |
| Sum of Squares | SS = 7343.75 |

Solution

$$\sigma = \sqrt{\frac{\sum_{i=1}^{n}(x_i - \mu)^2}{n}}$$

$$\sigma = \sqrt{\frac{SS}{n}}$$

$$\sigma = \sqrt{\frac{7343.75}{20}}$$

$$\sigma = \sqrt{\frac{7343.75}{20}}$$

$$\sigma = \sqrt{367.1875}$$

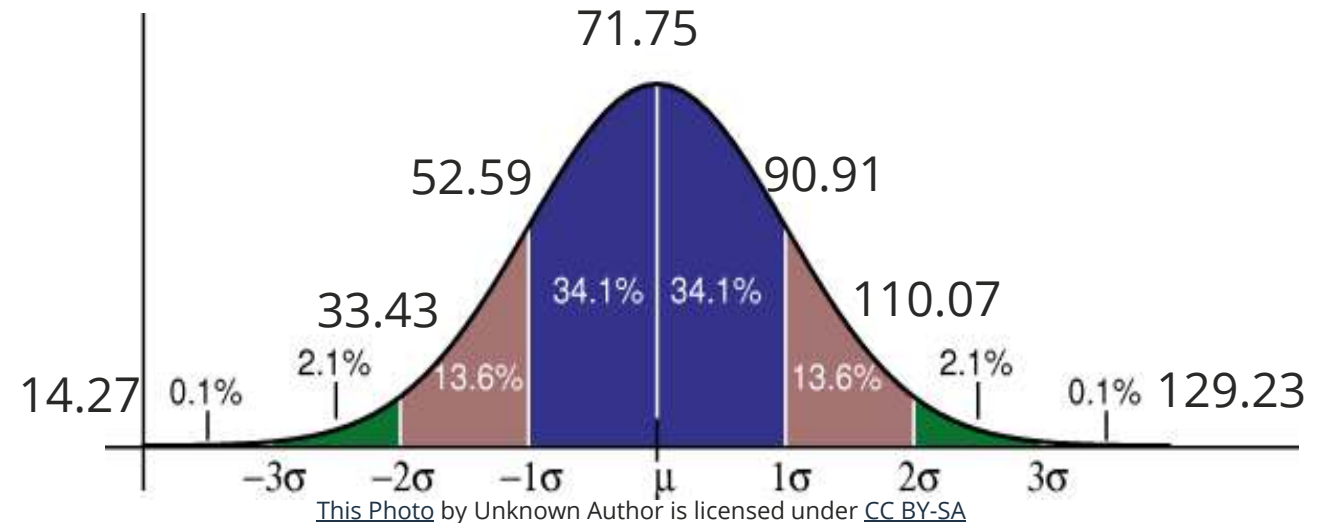$$\sigma = 19.162137$$



*This Photo by Unknown Author is licensed under CC BY-SA*

*https://www.calculatorsoup.com/calculators/statistics/standard-deviation-calculator.php*

# Plot of Mean = 71.75, σ = 19.16



**General Rule**
Percent of total data
68% – 95% – 99%
1σ       2σ       3σ

71.75

52.59          90.91

33.43          110.07

14.27          129.23

34.1%   34.1%

2.1%   13.6%        13.6%   2.1%

0.1%                                    0.1%

−3σ   −2σ   −1σ   μ   1σ   2σ   3σ

**So, out of sample of 20:**

**68% (13.6) fall between 52.59 and 90.91:** 30, 36, 45, 45, 62, 65, 67, 67, 76, 77, 77, 77, 79, 84, 85, 88, 90, 93, 95, 97
**95% (19.0) fall between 33.43 and 110.07:** 30, 36, 45, 45, 62, 65, 67, 67, 76, 77, 77, 77, 79, 84, 85, 88, 90, 93, 95, 97
**99% (19.8) fall between 14.27 and 129.23:** 30, 36, 45, 45, 62, 65, 67, 67, 76, 77, 77, 77, 79, 84, 85, 88, 90, 93, 95, 97
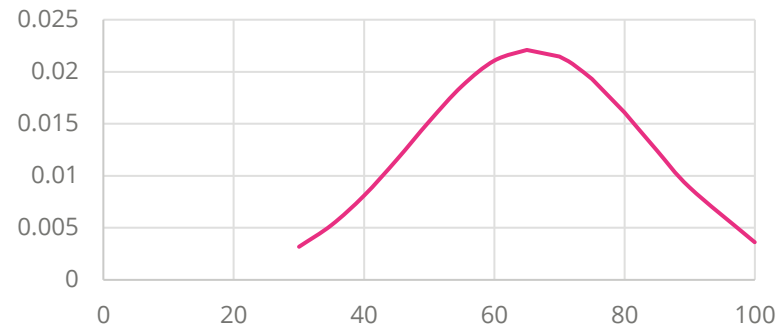
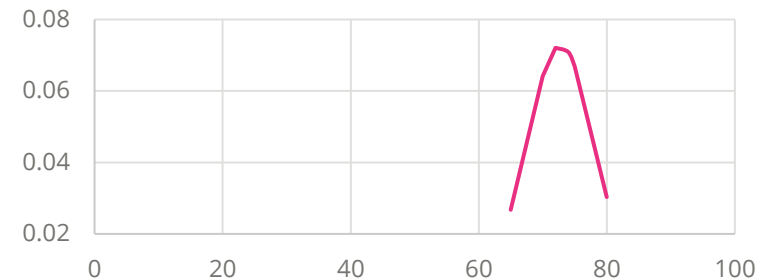# Acceptable Standard Deviation

- ## Norm-Referenced Test
  - 12.00 to 18.00



- ## Criterion-Referenced Test
  - No specific range
  - Should be small

Standard Deviation = 18
Mean = 65.6



Standard Deviation = 5.5
Mean = 72.75

# Variance

- A measure of how spread out a data set is – the average distance of a set of variables from the average value – used to calculate the standard deviation

- The more spread the data – the larger the variance in relation to the mean

    o 30, 45, 67, 90, 77, 36, 97, 65, 84, 79, 77, 62, 88, 67, 77, 95, 45, 85, 76, 93

    - This data set has a large spread – low is 30, high is 97
    - Variance = 367.1875
    - $\sqrt{367.1875} = 19.16$ which is the SD

| Variance | $\sigma^2 =$ 367.1875 |
|---|---|
| Standard Deviation | $\sigma =$ 19.162137 |
| Count | $n =$ 20 |
| Mean | $\mu =$ 71.75 |
| Sum of Squares | SS = 7343.75 |

Solution

$$\sigma^2 = \frac{\sum_{i=1}^{n}(x_i - \mu)^2}{n}$$

$$\sigma^2 = \frac{SS}{n}$$

$$\sigma^2 = \frac{7343.75}{20}$$

$$\sigma^2 = \frac{7343.75}{20}$$

$$\sigma^2 = 367.1875$$

## So What?!

*https://www.calculatorsoup.com/calculators/statistics/variance-calculator.php*

- Smaller variance means the SD is smaller which means the scores are closer together

    o **CRT** – closer is better – shows mastery of the test takers as a whole – **Smaller variance**
    o **NRT** – spread out provides better discrimination for rank order – **Larger variance**

# Standard Error of the Mean (SEM)

- Measures how precise the mean of a sample is as an estimate of the true mean of a population

- SEM = SD / √ of sample size
  - SD = 19.16
  - Sample size = 20
  - √ 20 = 4.472
  - 19.16 / 4.472 = 4.284

Theoretically, this means the likely error of a whole population based on the sample of 20 is ±4.284

**So What?!**

- The **smaller the SEM**, the more precise or closer to the predicted population mean our test is – **test is doing well!**

# Standard Error of Measurement (SE*m*)

- The standard error of measurement (SE*m*) is a measure of how much measured test scores are spread around a "true" score

- The SE*m* is especially meaningful to a test taker because it applies to a single score and it uses the same units as the test

- Calculated using the reliability value of the test

**So What?!**

- The **smaller** the SE*m*s the **greater precision** in the estimation of test-taker achievement

- The **larger** the SE*m*, the **less sensitive** our ability to detect changes in student achievement

# Skew

- How data looks when plotted

  - **Negative** values usually indicate a relatively **easy** test.  A negative skewness is said to be "skewed left" which means the left "tail" is longer relative to the right "tail"

  - **Positive** values usually indicate a **difficult** test.  A positive value is "skewed right" in that the right tail is longer relative to the left tail

  - Acceptable range -3 to +3



Negative Skew



Positive Skew

# Kurtosis

- Kurtosis is a measure of "peakness" of a distribution.  Another way to view this is flatness opposed to pointed when compared to a "normal" distribution curve.
  - A "**normal**" kurtosis, which is very rare, will have a value of **0.00**
  - **High** kurtosis value indicates a distinct peak near the mean that declines rapidly and has a heavy tail - common in a **CRT** but not desired in an NRT
  - **Low** kurtosis value indicates a relatively flat top near the mean - desired in an **NRT**
- Acceptable range -10 to +10

**High Kurtosis**

**Low Kurtosis**

# Reliability (Cronbach's Alpha)

- The extent to which the test or test-item is effective at measuring anything at all consistently

| Interpretation of Reliability Coefficient | |
|---|---|
| **Reliability** | **Interpretation** |
| .90 and above | Excellent reliability |
| .80 - .89 | Very good |
| .70 - .79 | Good – in range of most. Some items could be improved. |
| .60 - .69 | Somewhat low. Tests should be supplemented by other measures to determine grade or performance. Probably some items need improvement. |
| .50 - .59 | Suggests need for revision of test unless it is very short (ten or fewer items). Test definitely needs to be supplemented by other methods to determine grade. |
| Below .50 | Questionable reliability. The test should not contribute to the candidates grade and needs revision. |

*Adapted from SCOREPAC® Item Analysis, https://www.washington.edu/assessment/scanning-scoring__trashed/scoring/reports/item-analysis/*

# So...Now What Do We Know About the Test?

**41 is the value at the midpoint – ½ above and below**

**The average test score**

**Smaller SE*m* = Greater precision**

**Small SEM – Test is doing well!**

## Table of Test Statistics

| | | | | | |
|---|---|---|---|---|---|
| Number of examinees | 27 | Mean | 38.96/64 (60.88%) | Standard error of mean | 1.77/64 (2.77%) |
| Number of items | 40 | Median | 41/64 (64.06%) | Standard error of measurement | 4.13/64 (6.45%) |
| Maximum possible score | 64 | Mode | 45/64 (70.31%) | Skew | -3.038 |
| Minimum achieved score | 0/64 (0%) | Standard deviation | 9.21/64 (14.39%) | Kurtosis | 12.444 |
| Maximum achieved score | 51/64 (79.69%) | Variance | 84.88/64 (132.62%) | Test reliability (Cronbach's Alpha) | 0.799 |

*Reliability is most meaningful if all items cover the same subject area.*

**45 appears most often**

**Negative Skew means easy test**

**The data has a slight spread: OK for CRT**

**Function of the SD: √84.88 = 9.21**

**High kurtosis = scores clustered together**

**Test has good to very good reliability**

# Commonly Reported Test-Item Statistics

How is each test-item doing?

# Item Analysis Report (*Questionmark Analytics Reports*)

| Perception question id | 0000100001641014 | | |
|---|---|---|---|
| Question type | Multiple Choice | Question status | Normal |
| Question minimum possible score | 0 | Question maximum possible score | 1 |
| Number of participants presented the question | 28 | Number of participants who responded to the question | 25 |
| Item difficulty p-value | ◆ 0.52 (+/− 0.1) | Item reliability | 0.173 |
| Item-total correlation discrimination | ▤ 0.339 (−0.2/+0.173) | Item-rest correlation discrimination | ▤ 0.247 (−0.208/+0.187) |
| High-Low discrimination | 0.50 | | |
| Participant comments | *No comments were entered for this question* | | |

| Answer option information | | Number and percentage of participants achieving scores | | | |
|---|---|---|---|---|---|
| Outcome # | Answer option | All | Upper 27% | Middle 46% | Lower 27% |
| ✅ 1 | Item-total outcome correlation is a point-biserial calculation that compares a test item's score with the test taker's total exam score | 13 (46.4%) | 6 (75%) | 5 (41.7%) | 2 (25%) |
| 2 | Item-total outcome correlation is the measure of the number of people who answered a specific item correctly | 6 (21.4%) | 2 (25%) | 2 (16.7%) | 2 (25%) |
| 3 | Item-total outcome correlation is a function of the bell curve that defines the average deviation or degree of distribution of scores from the mean score | 4 (14.3%) | 0 (0%) | 2 (16.7%) | 2 (25%) |
| 4 | Item-total outcome correlation is a measure of how well a test has the capacity to repeat the same statistical results repeatedly | 2 (7.1%) | 0 (0%) | 2 (16.7%) | 0 (0%) |
| 5 | No response | 3 (10.7%) | 0 (0%) | 1 (8.3%) | 2 (25%) |
| Total assessment mean score | | 58.7 % | 71.7 % | 62.9 % | 39.5 % |

THE TWO

P-value    d-value

MOST IMPORTANT

# P-Value

- The P-value is also called the difficulty index
  - Correct response P-val – percentage of test-takers responding correctly
  - Incorrect response P-val – percentage of test-takers responding to each distractor

- CRT – Correct response P-val should be .80 or higher

- NRT – Correct response P-val should be .28 to .80

- Incorrect response P-vals should be somewhat equal

# *d*-value

- *d*-value represents the statistical function, "Point-Biserial Correlation," commonly called the discrimination index

- It is the degree to which the test item differentiates between those who know the material well and those who do not

- **Correct response** *d*-value should always be a **positive** number and the **incorrect responses** *d*-value should always be **negative**

- The range of the *d*-value is negative (-) 1.000 to positive (+) 1.000

    o Good correct item *d*-value range is +0.250 to +0.750
    o Good incorrect item *d*-value range is -0.250 to -0.750

# So What?!

- A **large positive** $d$-value such as .40 for the correct answer means that test takers with high scores on the test are also getting the item correct and those with lower test scores are getting the item wrong

- A **low** $d$-value for the correct response implies that test takers who get the item correct tend to do poorly on the exam overall and those who do well on the exam tend to get the item wrong

- A **negative** $d$-value for a correct response indicates there is some deficiency in the item which may include: item keyed incorrectly, item poorly constructed, misleading distractors, content inadequately taught, etc.

# *d*-value Discrimination Range

| General Discrimination Range | |
|---|---|
| **Absolute Value Range** | **Item Quality** |
| 0.50 or higher | Very high discrimination |
| 0.30 to 0.49 | High discrimination – possible item revision |
| 0.16 to 0.29 | Moderate discrimination – item needs revision |
| 0.15 or less | Low discrimination – review item to determine reason – possibly remove item |

**Note:** Negative values are expected for incorrect responses.  If a correct response has a negative value a problem is indicated.

*Adapted from: Pope, G. 2009, Item analysis analytics. Questionmark Corporation. Retrieved January 17, 2013, from http://www.questionmark.com/us/whitepapers/index.aspx*

# P-value and *d*-Value Generalizations

Observe relationships between P and *d* values

- Very easy or very difficult items have very little discrimination – do not tend to separate test takers who fully understand material from those who don't

- Items of moderate difficulty (60% - 70% correct response P-value) are generally more discriminating

- If all test-takers get item correct (P=1.00) or incorrect (P=0.00) the item does not discriminate at all – should be considered for removal – especially on NRT

- If all test-takers respond incorrectly (P=0.00) on CRT item – it may be keyed wrong or taught inadequately

# Sample P-value and d-value Interpretations

| Item Number | Response "A" | | Response "B" | | Response "C" | | Response "D" | |
|---|---|---|---|---|---|---|---|---|
| | P | $d$ | P | $d$ | P | $d$ | P | $d$ |
| 1 | **0.396** | **0.261** | 0.271 | -0.207 | 0.208 | -0.112 | 0.125 | 0.031 |
| 2 | **0.396** | **0.370** | 0.000 | 0.000 | 0.083 | -0.186 | 0.521 | 0.629 |
| 3 | 0.208 | 0.260 | **0.208** | **0.030** | 0.063 | -0.207 | 0.521 | -0.135 |
| 4 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | **1.000** | **1.000** |
| 5 | 0.417 | 0.061 | 0.125 | -0.214 | 0.271 | 0.231 | **0.188** | **-0.159** |

Responses keyed as "correct" are indicated by **bold underline**

**Note:** This report is not available in Questionmark analytics

# Item-Total Outcome Correlation Discrimination

- Item-total outcome correlation is a point-biserial calculation that compares a test item's score with the test taker's total exam score

- Higher item scores should mean higher exam scores overall

- **High** item-total correlation represents a **higher internal test consistency** and reliability
  - It means that test takers that score high on the test also scored higher on the test item than test takers that scored low on the test

- **Low** item-total correlation means that the test takers who scored low on the test are getting the answer correct more often than the test takers who scored high on the test
  - **Item needs review** – it may be confusing test-takers who are more competent

# Upper and Lower 27%

- Total-item discrimination is calculated using statistics for the test takers who score in the upper 27% minus statistics for those who score in the lower 27%.

- There should be a large positive difference between low and high

| LOWER 27% | MIDDLE 46% DISREGARDED | UPPER 27% |

# Interpretation of Total Outcome Correlation

| Total Outcome Correlation | Interpretation |
|---|---|
| Negative | Major problem indicated if this is occurring for a correct response – find out why |
| Around zero | No relationship between the test item score and the total assessment score – Review the items to determine why |
| 0 to 0.19 | Low correlation between outcome scores and assessment scores |
| 0.20 to 0.29 | Moderate correlation between outcome scores and assessment scores |
| 0.30 to 0.44 | Strong correlation between outcome scores and assessment scores |
| 0.45 or greater | Very strong correlation between outcome scores and assessment scores |

*Greg Pope – Item Analysis Analytics: The White Paper -https://www.questionmark.com/item-analysis-analytics-the-white-paper/*

# Item-Rest Correlation Discrimination

- Sometimes called Item-Remainder Correlation

- Most useful for short assessments of 25 items or fewer, small sample sizes, or assessments with different weighted items

- Good values typically 0.20 – 0.40 for cognitive tests and higher for typical-behavior tests

**So What?!**

- What's the difference between item-test (total) and item rest correlation?
  - Item-test correlation – shows how highly correlated the item is with the overall results
  - Item-rest correlation – shows the correlation of the item without including the item in the calculation (the rest of the items)

# High – Low Discrimination (D)

- Subtract the percentage of low-scoring participants who got the item correct from the percentage of high-scoring participants who got the item correct
  - **Example:** If 30% of low-scoring participants answered correctly, and 80% of high-scoring participants answered correctly, then the High-Low Discrimination is 0.80 – 0.30 = 0.50

**So What?!**

- **Positive** values indicate **good** discrimination, values **near zero** indicate that there is **little discrimination**, and **negative** discrimination indicates that the item is **easier** for low-scoring participants

| D Range | Interpretation |
|---------|----------------|
| 0.40 ≤ D ≤ 1.00 | Satisfactory High – Low discrimination |
| 0.30 ≤ D ≤ 0.40 | Some revisions may be required to the item |
| 0.20 ≤ D ≤ 0.30 | The item need revision |
| -1.00 ≤ D ≤ 0.20 | The item needs to be removed or completely revised |

*From: Measuring Educational Achievement – Robert L. Ebel, 1965*

# So...How Did This Test-Item Do?

| Perception question id | 0000100001641014 | | |
|---|---|---|---|
| Question type | Multiple Choice | Question status | Normal |
| Question minimum possible score | 0 | Question maximum possible score | 1 |
| Number of participants presented the question | 28 | Number of participants who responded to the question | 25 |
| Item difficulty p-value | ◆ 0.52 (+/- 0.1) | Item reliability | 0.173 |
| Item-total correlation discrimination | ▤ 0.339 (-0.2/+0.173) | Item-rest correlation discrimination | ▤ 0.247 (-0.208/+0.187) |
| High-Low discrimination | 0.50 | | |
| Participant comments | No comments were entered for this question | | |

| Answer option information | | | | All | Upper 27% | Middle 46% | Lower 27% |
|---|---|---|---|---|---|---|---|
| Outcome | | r option | test item's score with the test taker's total exam score | | | | |
| ✓ 1 | | otal outcome | relation is a poi | 13 (46.4%) | 6 (75%) | 5 (41.7%) | 2 (25%) |
| 2 | | otal outcom | relation is the n ... nswered a specific item correctly | 6 (21.4%) | 2 (25%) | 2 (16.7%) | 2 (25%) |
| 3 | Item-total outc | relation is a function of the bell curve that defines the average deviation or degree of distribution of scores fr... the mean score | 4 (14.3%) | 0 (0%) | 2 (16.7%) | 2 (25%) |
| 4 | Item-total out | relation is a measure of how well a test has the capacity to repeat the sa... ...edly | 2 (7.1%) | 0 (0%) | 2 (16.7%) | 0 (0%) |
| 5 | No resp | | 3 (10.7%) | 0 (0%) | 1 (8.3%) | 2 (25%) |
| Total assessment mean score | | | | 58.7 % | 71.7 % | 62.9 % | 39.5 % |

Callout annotations:
- Item reliability is questionable
- Item-rest correlation is OK but could be stronger – item may not correlate with other items
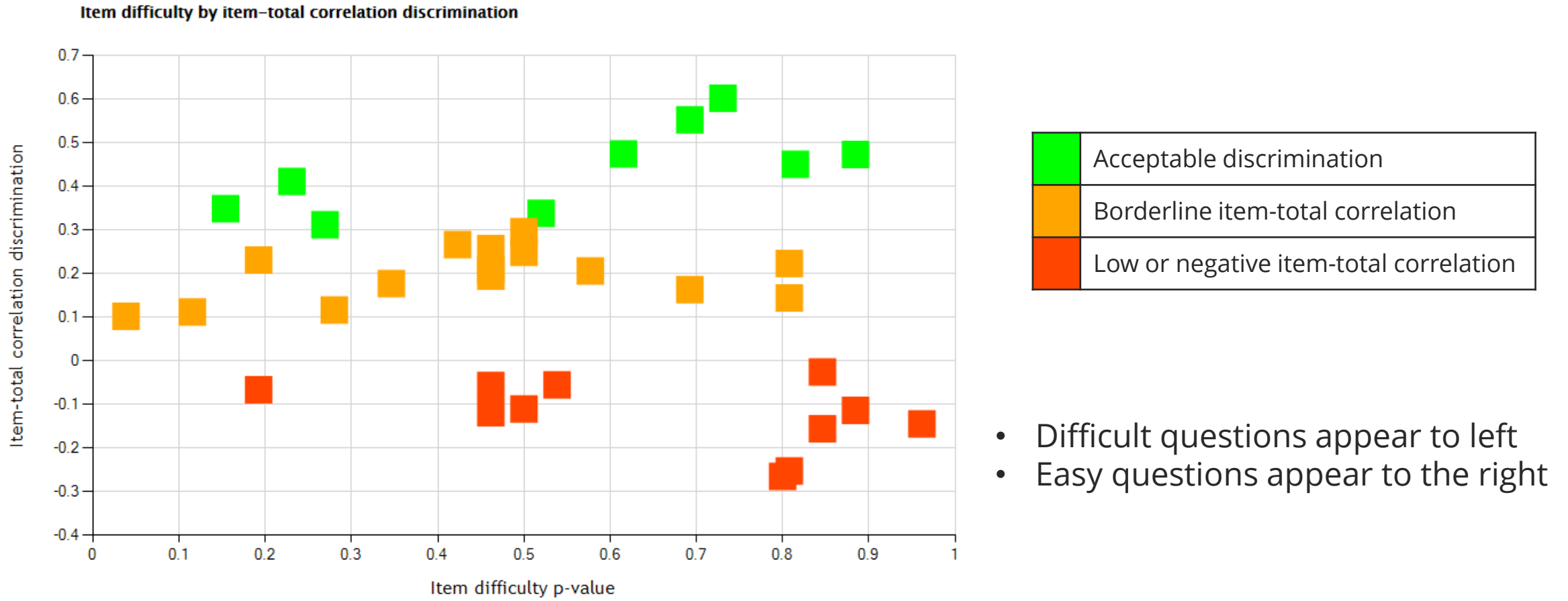- 52% are answering correctly – good value
- Item has good correlation with other items – could be better
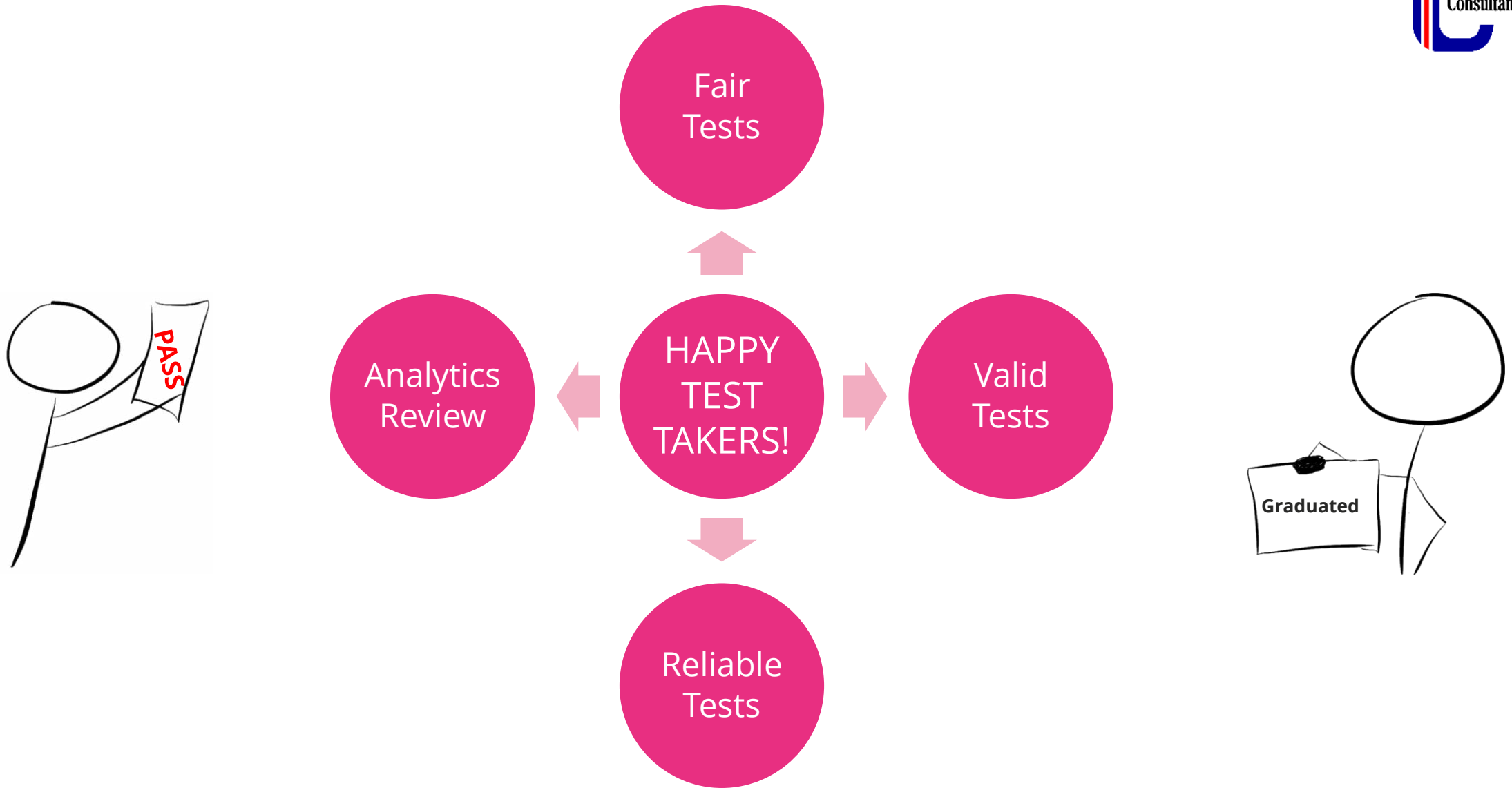- Item discriminates well – Right test-takers get it correct
- P-values distributed OK

# Example P-value vs. *d*-Value from Questionmark Analytics

**Item difficulty by item–total correlation discrimination**



Polytomous items are not plotted on this chart.

| | |
|---|---|
| 🟩 | Acceptable discrimination |
| 🟧 | Borderline item-total correlation |
| 🟥 | Low or negative item-total correlation |

- Difficult questions appear to left
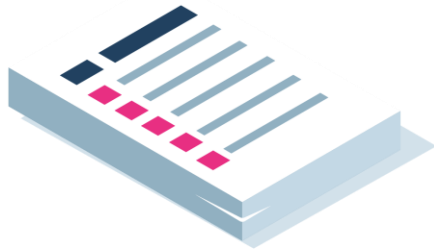- Easy questions appear to the right

# Questions?

question mark

# White papers, infographics, reports, eBooks and more!

## VIEW NOW:

Psychometrics – A collection of articles from Questionmark:
https://www.questionmark.com/category/psychometrics/

Webinar - Psychometrics 101: What your Psychometrician is REALLY saying?

# Upcoming webinars

## Introduction to Questionmark's Assessment Platform

◆ June 8, 2021 - 12:00 pm to 1:00 pm (EDT)

Learn the basics of authoring, delivering and reporting on surveys, quizzes, tests and exams. This introductory webinar explains and demonstrates key Questionmark features and functions.

**Click to Register**

## Beyond Recall: Taking Competency Assessments to the Next Level

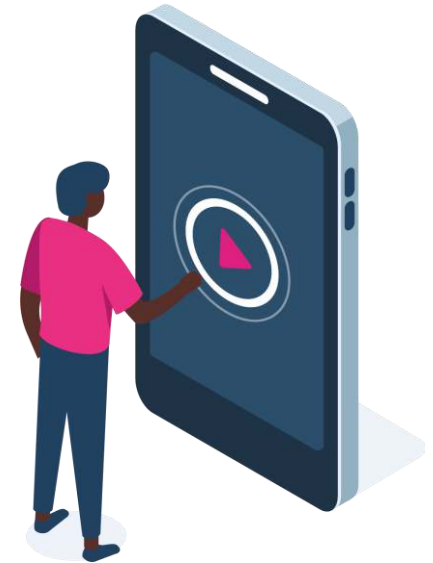◆ June 22, 2021 - 11:00 am to 12:00 pm (EDT)

Is it possible to assess someone's abilities to make judgments and decisions when they are faced with a dilemma? This session gives a general overview of why it's important to go beyond recall in competency assessments.

**Click to Register**

## Introduction to Questionmark's Assessment Platform

◆ June 24, 2021 - 10:00 am to 11:00 am (EDT)

Learn the basics of authoring, delivering and reporting on surveys, quizzes, tests and exams. This introductory webinar explains and demonstrates key Questionmark features and functions.

**Click to Register**

# questionmark

# Thank you for attending!

*Reach out to Questionmark at [sales@questionmark.com](mailto:sales@questionmark.com)
or request a demo at [https://www.questionmark.com/request-demo](https://www.questionmark.com/request-demo)*

*If you would like to reach out to Jim Parry – [james.parry@gocompassconsultants.com](mailto:james.parry@gocompassconsultants.com)
[www.gocompassconsultants.com](http://www.gocompassconsultants.com)*