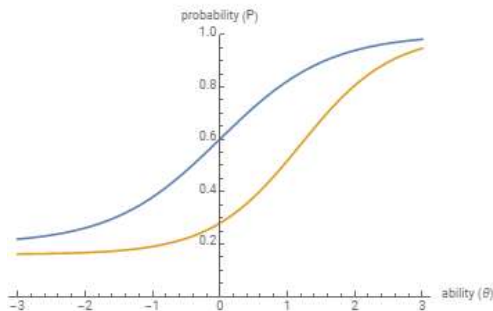


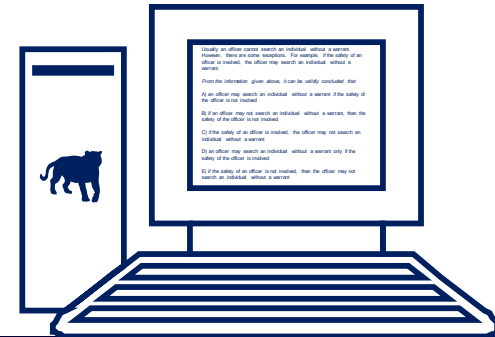
# Tutorial: How to Develop and Implement Unidimensional Computer Adaptive Tests

## A Tutorial



Courtesy Wolfram Demonstrations Projects

...mndtirioqlA...loop metl...noitslumia  
2, 3-PLM...θ, a, b, c...Monte Carlo



Scott K. Burtnick (U.S. Customs and Border Protection)  
Jeffrey M. Cucina (U.S. Customs and Border Protection)  
Kevin A. Byle (U.S. Customs and Border Protection)

The views expressed in this paper are those of the authors and do not necessarily reflect the views of U.S. Customs and Border Protection or the U.S. Federal Government.

# Who we are: US Customs and Border Protection



- US Customs and Border Protection (CBP)
  - America's unified border agency
  - Twin goals of anti-terrorism and facilitating legitimate trade and travel
  - Secures 328 ports of entry into US and borders between ports
  - Prevents narcotics, agricultural pests, smuggled goods and inadmissible visitors (e.g., aliens with outstanding criminal warrants) from entering US
  
- Personnel Research and Assessment Division (PRAD)
  - Part of CBP's Office of Human Resources Management
  - Group of I/O psychologists who design, develop, validate, and implement wide range of competency-based assessments and conduct organizational development work (e.g., survey research)
    - Entry-level and promotional assessments



# Who we are: US Customs and Border Protection



...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo



## Scott K. Burtnick

- Personnel Research Psychologist
- Expertise in psychometrics, test development, item response theory, differential item functioning, and data analysis.
- Leads survey research and conducts psychometrics analyses for CAT programs at U.S. Customs and Border Protection
- Previous work experience as psychometric lead for DoD security and intelligence certifications at Global Skills X-Change and conducting research for the U.S. Merit Systems Protection Board



## Jeffrey M. Cucina

- Personnel Research Psychologist
- Expertise in psychometrics, test development, item response theory, data analysis, and individual differences research.
- Psychometric lead for CAT programs at U.S. Customs and Border Protection.
- Leads criterion-related validity studies, test development, and job analyses.



U.S. Customs and  
Border Protection

# Who we are: US Customs and Border Protection



## Kevin A. Byle

- Supervisory Personnel Research Psychologist
- Expertise in psychometrics, test development, item response theory, and data analysis.
- Leads job analyses, test development and validation, and survey research at U.S. Customs and Border Protection.





# Overview of Tutorial

This tutorial covers the development of unidimensional computer adaptive tests (CATs) for dichotomously scored items

1. Overview of CATs
2. Brief review of item response theory (IRT)
3. Item pool development
4. Monte Carlo simulations for selecting algorithms
5. Scaling/equating/metric issue
6. Experimental item collection strategies
7. Creating instructions for programmers
8. User testing
9. Implementing CATs
10. Lessons learned/things to consider





1...mlitroglA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Overview of CAT



- **Advantages vs. Static Testing**

- **Testing Time:**

- Has the potential to reduce testing times by estimating ability levels with fewer items than static tests.

- **Test Security:**

- Offers more test security because items are drawn from very large test banks, and each test taker receives a different set of items.
    - This is becoming increasingly important as organizations increasingly are shifting to more mobile or remote forms of employment testing.







Monte Carlo  
a, b, c...  
2, 3-PLM...  
...noitlumj2  
...lo09  
...A...  
...mitropA...

# Overview of IRT



- Item Response Theory (IRT) uses item parameters to evaluate items on an Item Characteristic Curve (ICC).
  - $a$  (item discrimination)
  - $b$  (item difficulty)
  - $c$  (pseudoguessing)
- IRT relies on several assumptions that must be tested before computing IRT item parameters and developing CATs.
  - Dimensionality (one vs. many)
  - Examinee independence
  - Test item independence



1...mhtropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Overview of IRT



- Several considerations should be made when reviewing the ICCs including whether:
  - The set of items contains an adequate distribution of item difficulties.
  - Item discrimination meets key thresholds.
  - The pseudoguessing parameter is appropriate.
- Examining ICCs will help you decide whether to retain or discard items from the item bank, as well as whether a two or three parameter logistic model will best fit the data.



1...mhtropA...loo9 mell...noitlumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- In order to develop a CAT, an item pool of “hundreds, and possibly thousands” of items is needed (Hambleton, Swaminathan, & Rogers, 1991, p. 149).
- Data for each item would need to be collected on large samples:
  - 1-PLM (100-500 examinees; Wright, 1977; de Ayala, 2009, p. 42)
  - 2-PLM (500 examinees with 20 or more items; de Ayala, 2009, p. 105)
  - 3-PLM (BILOG works well with 1,000 examinees and 20 items [Mislevy, 1986]; 1,000+ examinees “strongly recommended” [de Ayala, 2009, p. 131])



...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Of course, the rules-of-thumb on the previous page are guidelines.
  - What is needed for your CAT depends on your goals (e.g., unproctored CAT would likely need a large number of items).
  - Regardless, you need to collect data on a lot of items using a lot of examinees (many more than for a static test).
- Also need to consider that only a percentage of items that will be developed will survive content reviews and item analyses.
  - We suggest looking at past performance of your item writers for the particular content domain and assuming that an extra 10% of items might need to be dropped solely for IRT calibration reasons.



...mlitrogA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Let's suppose you want to develop an unproctored version of a 30-item cognitive static test.
  - Using 3-PLM with 5 item pools of 150 items each.
  - Assuming 75% of items would be retained after analyses, a total of 1,000 ( $150 \times 5 \div .75$ ) items need to be administered to 1,000 examinees each.



1...multitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Source of examinees

1. Operational examinees: Collect data on experimental items using the individuals who are operationally taking the test (e.g., job applicants, certification test examinees).

- Pros: This could be the best approach psychometrically, especially if the current pool of operational examinees matches those of future pools.
- Cons:
  - A potential issue might occur if the number of operational examinees is small.
    - Suppose there are 1,000 operational examinees per year and 20 experimental items per examinee – it would take 10 years to collect data for one CAT.
  - Might require increasing testing time or holding scores (more on this later).



...mlitropA...loo9 mell...noitlumiz  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Item Pool Development



- Source of examinees
  
- 2. Convenience sample with incumbents: Collect data using people who are currently in the position (e.g., current onboard employees, established certified professionals).
  - Pros: Might be able to collect data from a large population, without being dependent on examinee volume.
  
  - Cons:
    - Could be problematic if range restriction is present. Need good representation at the low end of ability to estimate  $c$  parameter.
  
    - Incumbents might have lower motivation than operational examinees.
      - Response rate might be low.
  
      - Effort put into responding might be lower than examinees.
  
    - Might incur significant salary costs and consume staff time.



...mlitropA...loo9 mell...noislumig  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Source of examinees
  
- 3. Convenience sample with non-incumbents: Collect data using people who are not currently in the position (e.g., use incumbents in similar occupation, undergraduates, Mechanical Turk).

  - Pros: Might be able to collect data from a large population very quickly with little cost.
  - Cons:
    - The convenience sample might not be reflective of the operational examinees
      - Differences in ability distributions, item responses, test-taking motivation, etc. could be problematic.
      - Manifestations of the target construct might also differ (e.g., if CAT will measure specialized knowledge, then couldn't collect data from the general population).





1...multitropA...loo9 mell...noislumig  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
  1. Common item/non-equivalent groups: Administer an anchor test along with different batches of experimental items. Rotate through different batches until 1,000 examinees have taken each item.
    - The anchor items could be the operational test.
    - Could also shorten the operational test and use experimental items to get full length
      - This would require holding scores and conducting item analyses and equating for each batch of items.
    - Need to consider examinee test-taking fatigue (although it may be no different than typical testing situations involving common item/non-equivalent groups equating).



1...multitropA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
  1. Common item/non-equivalent groups: Administer an anchor test along with different batches of experimental items. Rotate through different batches until 1,000 examinees have taken each item.
    - From an IRT perspective, you would conduct separate calibrations on each combined anchor/experimental test. Next, you would apply a scale transformation to place the theta, a, and b parameters on the same scale (more on this later).
    - Hambleton, Swaminathan, and Rogers (1991, p. 129) mention this is typically the most feasible design for IRT purposes.



...mltltropA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
- 2. Spiraling: Randomly assign examinees to take either operational test or a form containing only experimental items.
  - Would need to hold scores and conduct item analysis and equating (to place experimental items on scale of operational test) each time a new form of experimental items is introduced.
  - Wouldn't increase test-taking fatigue and could finish data collection quicker; however, it requires more item analysis and equating work to produce scores for operational purposes.
    - For each administration, would have to conduct item analysis and equating to produce operational scores for the experimental items.



...mlitropA...loo9 mell...noitlumj2  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
- 2. Spiraling: Randomly assign examinees to take either operational test or a form containing only experimental items.
  - From an IRT perspective, when spiraling is used, the item parameters and thetas for the two groups within one administration should be on the same scale (due to the random assignment; Kolen & Brennan, 2004, p. 166). However, this only works for the first administration when collecting experimental CAT items.
  - In subsequent administrations, would need to treat item parameters for operational (or already IRT-analyzed) test as fixed and conduct transformation of thetas and a and b parameters for new test to old test's scale.



1...multitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
- 3. Single group: One single group of examinees takes operational items and all experimental items. Could counterbalance order of operational and experimental items if desired.
  - We mention this approach for completeness. Unless the item pool will be very small or the items can be completed very quickly with little examinee fatigue, this approach is not going to be practical.
  - From an IRT perspective, only one calibration would be conducted using both the operational and experimental items. No scale transformation would be needed.
  - Common persons equating is a variant of this method. A single group takes both sets of items and two separate groups take either experimental or operational items.



1...multitrogA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
4. Harvesting old data: Use data from established forms to create an item bank.
- This approach assumes you have data for enough different forms of the test to create the CAT item pool. This might be possible for a large-scale testing program with a large number of equated forms and examinees.
  - Equating strategies 1-3 listed on the previous slides would have to have been implemented when the data were originally collected.
  - Would conduct IRT analyses on the existing data.
  - Need to consider scale drift and changes to the examinee population over time (e.g., item parameters for data collected 20 years ago might not be appropriate today).



1...multitropA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Data collection strategy
  
- 5. Combined static and CAT data collection: First collect data (using strategies 1-4) for experimental items for one CAT item pool using a static test. Next, develop the CAT and collect data for subsequent CATs by administering static experimental items alongside the CAT.
  - This approach assumes you are developing multiple CATs for a single test. Rather than collecting data on enough experimental items to develop multiple CATs, you first collect data for one CAT and implement it. After it is implemented, you collect data for the subsequent CATs.
  - From an IRT perspective, the pre-CAT data collection and calibration is the same as for strategies 1-4. After the CAT is implemented, a CAT pretest item calibration approach is used to place the experimental item parameters on the same scale as the operational CAT (more on this later).



Monte Carlo  
2, 3-PLM... $\theta$ , a, b, c...

# Item Pool Development



- Example for Strategy 1 (Common item/non-equivalent groups): Placing parameters from different calibrations onto same scale.
- Suppose you collected data for 20 batches of 10 experimental items and had 30 common items across all administrations.
- Each batch of 10 experimental items and 30 common anchor items would be calibrated separately.
- This yields 20 different IRT calibrations. Note that the scales for each calibration are different. For example, the mean  $\theta$  will be 0 in all calibrations, despite the fact that the examinees might have differed in ability levels across the different administrations.
  - This is because by default, BILOG and most other IRT programs assign a mean  $\theta$  value of 0 for each calibration.





# Item Pool Development



- You need to place the parameters from the 20 different IRT calibrations onto the same scale using the 30 common items as anchors.
  - You will want to choose 1 of the 20 calibrations to be the base form/calibration and transform the item parameters from each of the 19 remaining calibrations to those on the base.
- The below equations can be used to place the a and b parameters for two different tests (X and Y) on a common scale. The equations use the a and b parameters obtained from the separate calibrations of the common items (denoted c):

$$b_{Yc} = \alpha b_{Xc} + \beta$$

$$a_{Yc} = \frac{a_{Xc}}{\alpha}$$

(Hambleton, Swaminathan, & Rogers, 1991, p. 129)





# Item Pool Development

- Before we can use the equation, we need to estimate  $\alpha$  and  $\beta$ . There are a few different methods:
- Regression: Create regression equation to predict item parameters for test Y using those for test X. There is a symmetry issue since these will not be the same as for predicting X using Y; therefore, this method is not widely used.
- Mean/Sigma: Examines differences in means and SDs of anchor item parameters to compute  $\alpha$  and  $\beta$ :

$$\alpha = \frac{SD_{Yc}}{SD_{Xc}}$$

$$\beta = \overline{b_{Yc}} - \alpha \overline{b_{Xc}}$$





# Item Pool Development

- Mean/Mean: Used in 1-PLM Rasch Model:

$$\beta = \overline{b_{Yc}} - \overline{b_{Xc}}$$

- Robust Mean/Sigma: Similar to the Mean/Sigma method; however, it takes into account standard errors in estimating item parameters using a 5-step process (Linn et al., 1981). Stocking and Lord's (1991) adaptation takes into account outliers.
- Test Characteristic Curve: Iterative process that incorporates differences in discrimination (a) parameters (other methods focus on b) using test characteristic curve.
  - There are two variations of this: Haebara (1980) and Stocking and Lord (1983)
  - Overall, the Stocking and Lord (1983) method seems to be the preferred method in the literature (de Ayala, 2009; Kaskowitz & de Ayala, 2001; Baker & Al-Karni, 1991).



Monte Carlo  
2, 3-PLM... $\theta$ , a, b, c...

# Item Pool Development



- There are a number of programs that can be used to implement these methods:
  - EQUATE (Baker, Al-Karni, & Al-Dosary, 1991) and EQUATE 2.0 (1993).
  - equateIRT R package (Battauz, 2018; Wiberg, 2018)
  - We prefer IRTEq (Han, 2009), a free windows-based program that has a GUI and syntax option:

<https://www.umass.edu/remp/software/simcata/irteq/>

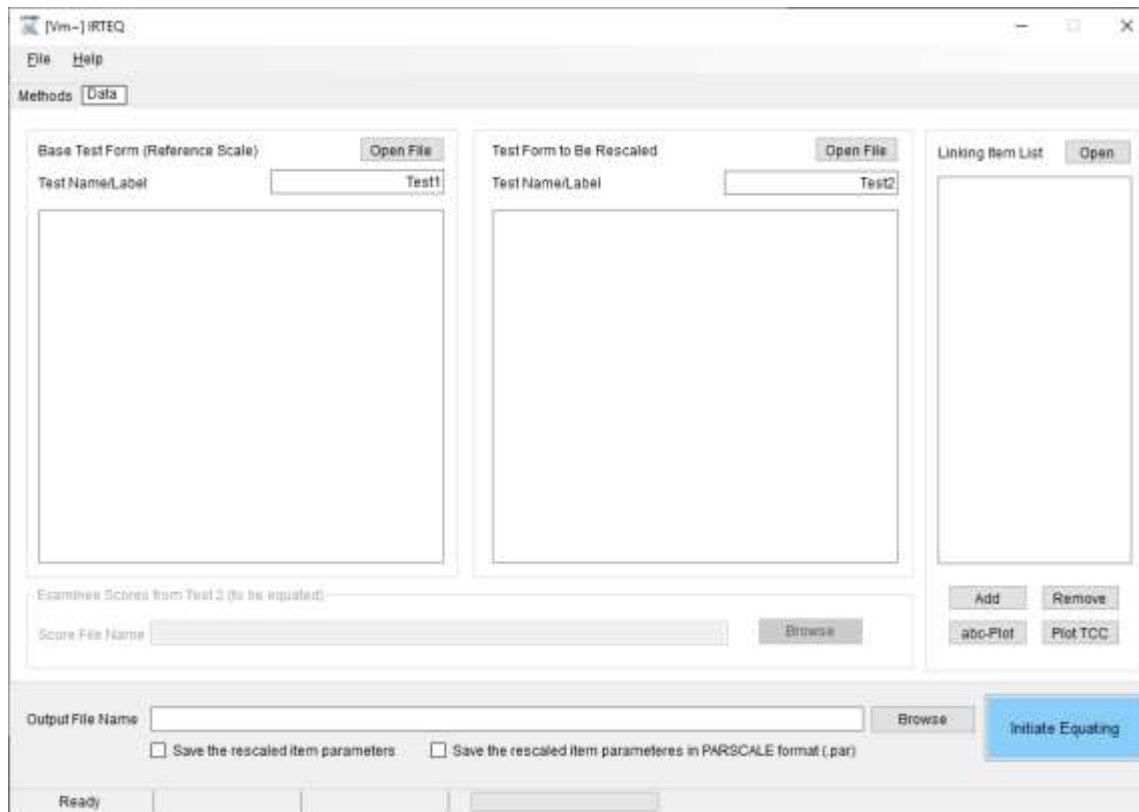


1...mltrpA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

## Item Pool Development



- Using IRTEq
  - First open IRTEq and navigate to the Data tab:



...mlitrogA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- IRTeq comes with several example files:
  - examinee\_700.wge (thetas for 700 examinees using a WinGen format)
  - example1.LIL (Linking Item Link text file showing which items in the two test forms are anchor items)
  - example1.syn (syntax file for the example)
  - test1.PAR (test 1 item parameters and standard errors in PARSCALE format)
  - test2.PAR (test 2 item parameters and standard errors)
  
- Depending on which IRT software you are using, you may need to manually convert your theta and item parameter output to match the ones that IRTeq accepts (i.e., WinGen and PARSCALE).



...mlitropA...loo9 melli...noislumia  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Let's begin by opening the data files for the two test administrations. Click on Open File for the Base Test Form and select test1.PAR and then click on Open File for Test Form to be Rescaled and select test2.PAR.

The screenshot shows the IRTEQ software interface with two data tables. The first table is titled 'Base Test Form (Reference Scale)' and the second is 'Test Form to Be Rescaled'. Both tables have columns for Test Name/Label, Item ID, and various statistical parameters.

Test Name/Label	Item ID	Item Type	Item Difficulty	Item Information	Item Discrimination	
Test1	1	3PLM	2	0.450	0.504	0.109
Test1	2	3PLM	2	0.290	-1.351	0.240
Test1	3	3PLM	2	0.467	-1.220	0.184
Test1	4	3PLM	2	0.371	-0.786	0.209
Test1	5	3PLM	2	0.593	-0.265	0.256
Test1	6	3PLM	2	0.559	-0.337	0.246
Test1	7	3PLM	2	0.485	-0.878	0.212
Test1	8	3PLM	2	0.715	-1.441	0.229
Test1	9	3PLM	2	0.539	-0.517	0.260
Test1	10	3PLM	2	0.467	-0.894	0.170
Test1	11	3PLM	2	0.443	-1.624	0.242
Test1	12	3PLM	2	0.738	-1.602	0.206
Test1	13	3PLM	2	0.579	-1.045	0.170
Test1	14	3PLM	2	0.435	-1.836	0.175
Test1	15	3PLM	2	0.407	-1.720	0.000

Test Name/Label	Item ID	Item Type	Item Difficulty	Item Information	Item Discrimination	
Test2	1	3PLM	2	0.714	-1.965	0.187
Test2	2	3PLM	2	0.546	-3.369	0.228
Test2	3	3PLM	2	0.265	-1.884	0.215
Test2	4	3PLM	2	0.294	-1.088	0.194
Test2	5	3PLM	2	0.537	0.068	0.462
Test2	6	3PLM	2	0.580	-1.871	0.175
Test2	7	3PLM	2	0.603	-2.317	0.208
Test2	8	3PLM	2	0.271	1.790	0.373
Test2	9	3PLM	2	0.578	-0.658	0.174
Test2	10	3PLM	2	0.612	-1.629	0.210
Test2	11	3PLM	2	0.658	-0.116	0.235
Test2	12	3PLM	2	0.481	0.406	0.299
Test2	13	3PLM	2	0.595	0.292	0.140
Test2	14	3PLM	2	0.675	-0.115	0.247
Test2	15	3PLM	2	0.332	-0.625	0.175



...mlitrogA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Item Pool Development



- Next, we need to tell IRTEq which pairs of items are the anchors in both forms. For example, suppose item 1 in test 1 is the same as item 46 in test 2. You can click on the parameters for both items and then click add to begin creating a Linking Item List.

The screenshot shows the IRTEQ software interface with three main panels:

- Base Test Form (Reference Scale):** Test Name/Label: Test1. Contains a table of items 1-6 and 18-19.
- Test Form to Be Rescaled:** Test Name/Label: Test2. Contains a table of items 44-49 and 61-62.
- Linking Item List:** Contains the number 1,46.

Below the tables is a section for "Examinee Scores from Test 2 (to be equated)" with a "Score File Name" input field and a "Browse" button. On the right side, there are buttons for "Add", "Remove", "abc-Plot", and "Plot TCC".

Item ID	Form	Count	Parameter 1	Parameter 2	Parameter 3
1	3PLM	2	0.450	0.504	0.109
2	3PLM	2	0.290	-1.351	0.240
3	3PLM	2	0.467	-1.220	0.184
4	3PLM	2	0.371	-0.786	0.209
5	3PLM	2	0.593	-0.265	0.256
6	3PLM	2	0.550	-0.337	0.246
18	3PLM	2	0.316	-2.052	0.213
19	3PLM	2	0.342	-1.468	0.221

Item ID	Form	Count	Parameter 1	Parameter 2	Parameter 3
44	3PLM	2	0.657	-0.298	0.201
45	3PLM	2	0.732	-0.080	0.181
46	3PLM	2	0.498	0.615	0.151
47	3PLM	2	0.363	-0.943	0.281
48	3PLM	2	0.543	-1.054	0.213
49	3PLM	2	0.313	-0.815	0.243
61	3PLM	2	0.472	-0.549	0.357
62	3PLM	2	0.418	-1.552	0.000



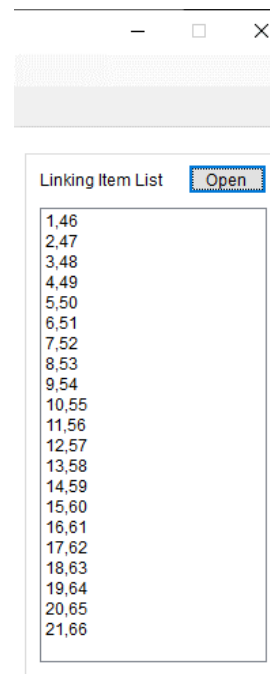
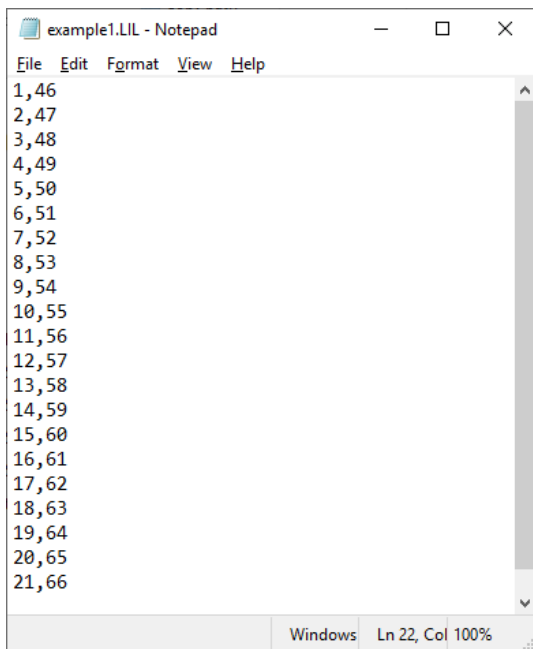


1...mhtropA...loo9 mell...notistumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

## Item Pool Development



- Or, you can create a text file that shows the item pairings using commas and import it into IRTeq using the Open button in the Linking Item List. The provided example1.LIL is an example of how to set this up. Make sure you save the file with a .LIL extension.

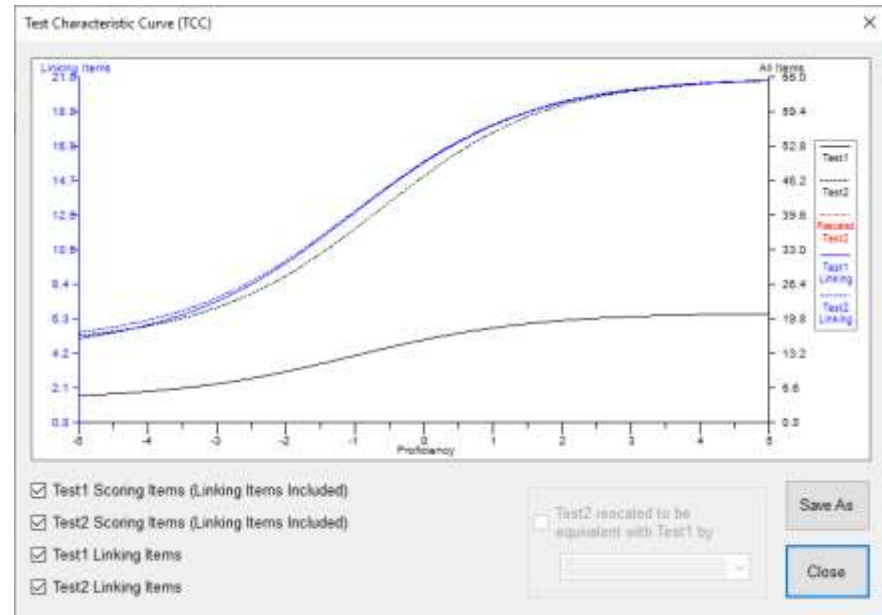
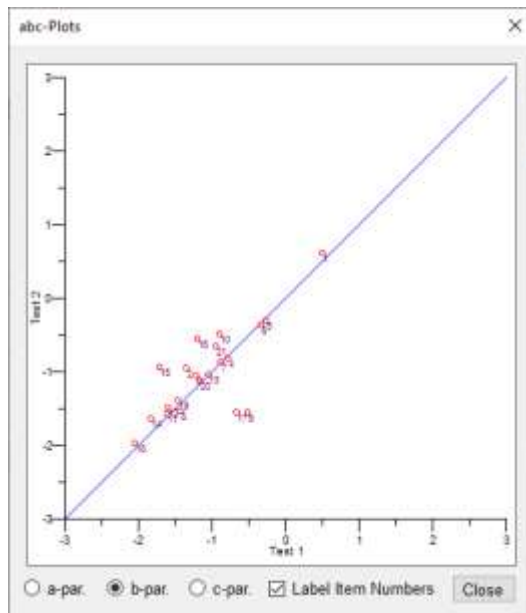


Monte Carlo  
2, 3-PLM...θ, a, b, c... Monte Carlo

## Item Pool Development



- (Side note: At this point, you can also click the abc-Plot and Plot TCC buttons in the Linking Item List area to explore plots of the item parameters across series for the anchor items).



1...mltltropA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Now, go back to the Methods tab and select the IRT equating methods you want to use under the Item Scaling Method section.
- For CAT development purposes, you can ignore the options for external linking items and true score equating.
- You can probably also ignore the option for averaging the parameters from the linking items.
  - Sometimes psychometrician equating only two forms will average the item parameters for the anchor items across the two calibrations. Since our example has 20 different calibrations (and IRTEq only handles 2 at a time), this approach won't work for us.
- Make sure that you use the same scale ( $D = 1$  or  $1.7$ ) that your IRT software used in your initial calibrations.



1...mlitropA...loo9 mell...noitlumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

## Item Pool Development



[Vm~] IRTEQ

File Help

Methods Data

Item Scaling Method

- Mean/Mean
- Mean/Sigma
- Robust Mean/Sigma
- TCC (Haebara)
- TCC (Stocking and Lord)

Linking Items

Internal (Scored)  External (Not Scored)

Choice of Weights for the TCC Scaling Methods

Uniform Distribution (Lowest:  Highest: )

Normal Distribution (Mean:  SD: )

Actual (examinee) Distribution

Scale

Logistic (D=1.0)  Normal Ogive (D=1.7)

True Score Equating

Save the Equated Scores

Produce the Conversion Table

Average the Rescaled and the Original Parameters for the Linking Items (an option with MMMS/RMS)

Directions

```
==== Scaling Item Parameters ====  
  
> For Robust Mean/Sigma Method, item parameter estimates and standard errors should be provided in PARSCALE parameter file format (*.par).  
  
> For Mean/Mean, Mean/Sigma, TCC(Haebara), TCC(Stocking & Lord) method(s), item parameter files are needed either in *.par (PARSCALE) or in *.wgi (WinGen).
```

Output File Name  Browse

Save the rescaled item parameters  Save the rescaled item parameters in PARSCALE format (.par)

Initiate Equating

Ready



1...mltrpA...loo9 mtl...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Item Pool Development



- Now we need to specify the name of an output file using the bottom left section of either the Methods or Data tab. You can also check the boxes for saving the rescaled item parameters.
- Then quick Initiate Equating on either tab to run the analysis.

15	3PLM	2	0.407	-1.720	0.000	58	3PLM	2	0.536	-1.025	0.233	17,62
16	3PLM	2	0.371	-1.202	0.238	59	3PLM	2	0.440	-1.626	0.215	18,63
17	3PLM	2	0.579	-0.665	0.230	60	3PLM	2	0.489	-0.940	0.233	19,64
18	3PLM	2	0.316	-2.052	0.213	61	3PLM	2	0.472	-0.549	0.357	20,65
19	3PLM	2	0.342	-1.468	0.221	62	3PLM	2	0.418	-1.552	0.000	21,66

Examinee Scores from Test 2 (to be equated)

Score File Name

Output File Name

Save the rescaled item parameters  Save the rescaled item parameters in PARSCALE format (.par)

Ready



Monte Carlo  
2, 3-PLM...θ, a, b, c...

# Item Pool Development



- An output file will be generated and should open automatically.
- It provides some summary info for the number of items, which you should verify.
- The output also gives the mean and SDs for the item parameters across the two calibrations and the correlations between these.

```

example1.OUT - Notepad
File Edit Format View Help
IRTEQ (Han, 2007)
output file (5/12/2020 4:17:23 PM)

=====

Test 1: C:\temp\IRTEQ\EXAMPLES\test1.PAR
# of items: 21

Test 2: C:\temp\IRTEQ\EXAMPLES\test2.PAR
# of items: 66

# of pairs of linking items: 21

=====

mean discrimination of linking items in Test1:      0.477
mean difficulty of linking items in Test1:         -1.072
standard deviation of difficulty parameters of linking items in Test1: 0.603

mean discrimination of linking items in Test2:      0.490
mean difficulty of linking items in Test2:         -1.007
standard deviation of difficulty parameters of linking items in Test2: 0.596

Correlation coeff. for discrimination of linking items: 0.824
Correlation coeff. for difficulty of linking items:    0.767
Correlation coeff. for guessing of linking items:    -0.056

=====
    
```



Monte Carlo  
2, 3-PLM... $\theta$ ,  $a$ ,  $b$ ,  $c$ ...

# Item Pool Development



- Once you scroll down, you'll see the  $\alpha$  and  $\beta$  coefficients (labeled as A and B) for each of the methods.
- The remainder of the output tells you where the rescaled parameter values are saved.

```

example1.OUT - Notepad
File Edit Format View Help
=====
Equating coefficient A with Mean-Mean method:      1.029
Equating coefficient B with Mean-Mean method:      -0.036

Equating coefficient A with Mean-Sigma method:     1.013
Equating coefficient B with Mean-Sigma method:     -0.052

Equating coefficient A with Robust Mean-Sigma method: 0.984
Equating coefficient B with Robust Mean-Sigma method: -0.130

Equating coefficient A with TCC (Haebara) method:   1.00
Equating coefficient B with TCC (Haebara) method:   0.03
Minimized Loss Function Value with TCC (Haebara) method: 0.00765

Equating coefficient A with TCC (Stocking & Lord) method: 0.99
Equating coefficient B with TCC (Stocking & Lord) method: 0.03
Minimized Loss Function value with TCC (Stocking & Lord) method: 0.00127

=====

Test2 items rescaled onto the scale of Test1 using Mean/Mean Method were saved in
example1.MM.wgi (5/12/2020 4:17:37 PM)
Test2 items rescaled onto the scale of Test1 using Mean/Sigma Method were saved
in example1.MS.wgi (5/12/2020 4:17:37 PM)
Test2 items rescaled onto the scale of Test1 using Robust Mean/Sigma Method were
saved in example1.RMS.wgi (5/12/2020 4:17:37 PM)
Test2 items rescaled onto the scale of Test1 using TCC(Haebara) Method were saved
Windows (CRLF) Ln 1, Col 1 100%

```







# Item Pool Development

...mlitropA...loo9 mell...notistumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

- In our example with 20 calibrations, you'll repeat this process 18 more times (1 for each of the remaining calibrations) and compile the  $\alpha$  and  $\beta$  coefficients. Next, you'll plug these into the equations shown earlier to place all of the item parameters on the same scale as the base series.
- IRTEq does have a syntax file option whereby you can use code (rather than pointing and clicking) to set up the analyses. Here is the syntax for the example:

```
example1.syn - Notepad
File Edit Format View Help
IRTEQ syntax file saved by a user
12/7/2007 12:55:03 PM
MM,MS,SL,SAVE,PARSCALE
-3,3,NORM,0,1,AVERAGE
LOGISTIC
INTERNAL
EQUATE,CONVERSION
C:\IRTEQ\EXAMPLES\example1.out
C:\IRTEQ\EXAMPLES\test1.PAR
C:\IRTEQ\EXAMPLES\test2.PAR
C:\IRTEQ\EXAMPLES\example1.lil
C:\IRTEQ\EXAMPLES\examinee_700.wge
Windows Ln 13, Co 100%
```

- You can also create a cue file that lists different syntax files and run a large number of syntax files automatically at once





1...multitopA...loo9 mell...noislumig  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- At this point, you should have a listing of all of the item parameters for your item pool on the same metric.
- If you have collected enough data for multiple CATs, you'll want to divide the items into different pools, trying to ensure that the test information functions follow the distributions of test scores (or are focused on the cut score) as much as possible.
- We will assume for now that you only have enough data for one CAT at this point.



# Monte Carlo Simulations



- We need to decide what type of algorithm to use. The choice is often guided by the characteristics of the item pool, your professional judgment, and the philosophy of the testing program.
  - At this point the literature doesn't explicitly recommend one algorithm over another.
    - The answer often depends on the characteristics of the item pool.
    - There are also too many possible algorithms to study all of them across different types of item pools.
    - So, you need to do a Monte Carlo simulation to select an algorithm (or confirm/compare ones you have in mind)



1...mltrpA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- A number of programs exist for conducting the Monte Carlo simulation:
  - CATSim (Weiss & Guyer, 2012): Available for purchase at <https://assess.com/catsim/>.
  - R packages: catIrt (Nydick, 2014); xxIRT (Luo, 2019)
  - We prefer SimulCAT (Han, 2012), a free windows-based program that has a GUI and syntax option:

<https://www.umass.edu/remp/software/simcata/simulcat/>



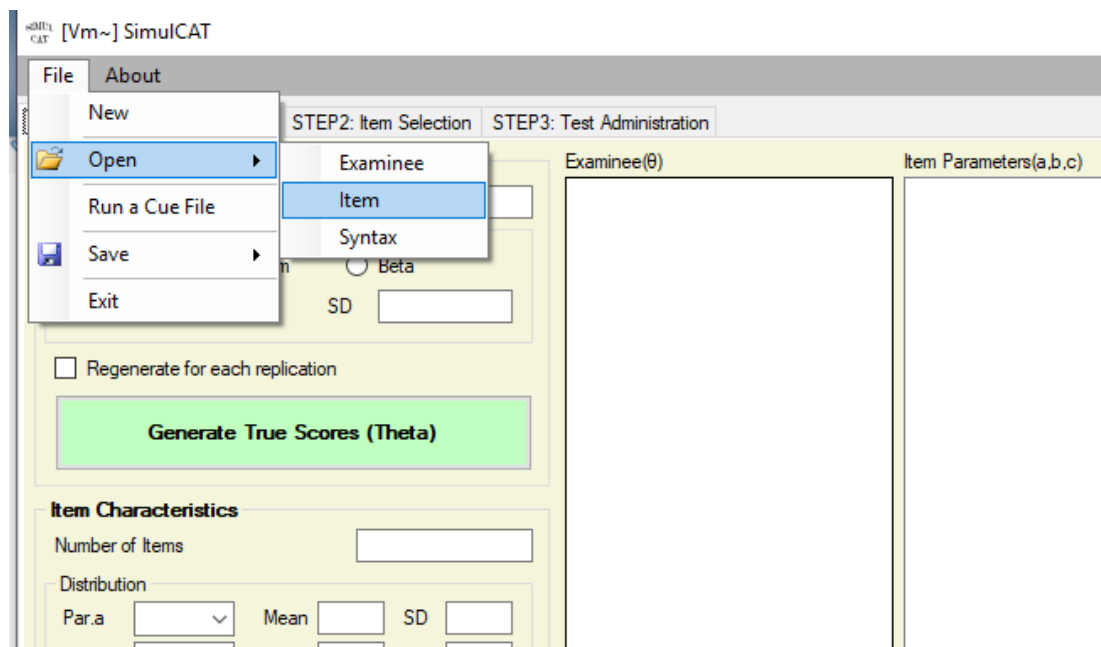
1...multitropA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations



- This program also comes with example files. We'll only be using one today: Example\_ItemPool500.wgix (this contains item parameters for 500 items using a 3-PLM).

- Let's open SimulCAT and load this file:





...mhtirpA...loo9 mell...noitlumj2  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations

- All of the item parameters will appear on the right side of the window:

- You can click plot items and view the ICCs, item information, TCC, and test information plots

The screenshot shows the SimuCAT software interface. The window title is "[Vm-] SimuCAT". The menu bar includes "File" and "About". The main area is divided into several sections:

- Examinee Characteristics:** Includes fields for "Number of Examinees", "Distribution" (Normal, Uniform, Beta), "Mean", and "SD". There is a checkbox for "Regenerate for each replication" and a green button labeled "Generate True Scores (Theta)".
- Item Characteristics:** Includes fields for "Number of Items", "Distribution", and three parameter sets (Para, Mean, SD). There is a checkbox for "Scale to normal metric (scaling factor 0=1.702)" and a checkbox for "Add to the previous item set". A red button labeled "Generate True Item Parameters" is at the bottom.
- Item Parameters(a,b,c):** A table listing 40 items with their parameters. The table has columns for Item ID, Distribution, and three parameters (a, b, c).
- Summary Statistics:** At the bottom, it shows "N = 0", "Mean = 0", "SD = 0", and "Histogram". To the right, it shows "a: Mean=0.800 / SD=0.222", "b: Mean=0.016 / SD=0.941", and "c: Mean=0.180 / SD=0.017". A "Plot Items" button is also present.

Item	Distribution	a	b	c	
1	3PLM	2	0.977	-0.553	0.157
2	3PLM	2	1.019	0.420	0.386
3	3PLM	2	0.670	-1.535	0.206
4	3PLM	2	0.911	-0.298	0.203
5	3PLM	2	0.872	0.091	0.180
6	3PLM	2	0.527	0.661	0.207
7	3PLM	2	0.699	2.507	0.176
8	3PLM	2	1.086	-0.258	0.153
9	3PLM	2	0.837	0.449	0.184
10	3PLM	2	1.115	-1.380	0.388
11	3PLM	2	0.918	-1.304	0.168
12	3PLM	2	0.624	1.132	0.185
13	3PLM	2	1.060	-1.123	0.185
14	3PLM	2	1.130	0.137	0.179
15	3PLM	2	0.966	0.326	0.206
16	3PLM	2	0.682	1.448	0.185
17	3PLM	2	0.645	0.330	0.187
18	3PLM	2	0.810	0.979	0.160
19	3PLM	2	0.550	0.201	0.181
20	3PLM	2	0.912	1.457	0.165
21	3PLM	2	0.771	1.476	0.155
22	3PLM	2	1.186	-1.467	0.589
23	3PLM	2	0.711	-0.430	0.151
24	3PLM	2	1.024	0.064	0.167
25	3PLM	2	0.857	0.907	0.170
26	3PLM	2	0.801	0.663	0.156
27	3PLM	2	0.485	0.219	0.156
28	3PLM	2	0.773	0.579	0.199
29	3PLM	2	0.519	-0.534	0.174
30	3PLM	2	0.622	0.797	0.185
31	3PLM	2	0.963	1.420	0.181
32	3PLM	2	0.964	0.034	0.154
33	3PLM	2	0.939	-0.606	0.203
34	3PLM	2	1.149	0.006	0.150
35	3PLM	2	0.757	-1.253	0.152
36	3PLM	2	0.671	0.715	0.168
37	3PLM	2	0.826	0.672	0.159
38	3PLM	2	0.930	-1.546	0.160
39	3PLM	2	0.460	-1.230	0.183
40	3PLM	2	1.046	-0.650	0.200





# Monte Carlo Simulations

- Next, you should specify the number of simulees/examinees (10,000 is usually a good number for a Monte Carlo simulation) and the mean and SD of the theta distribution (0,1).

The screenshot shows the SimulCAT software interface. The title bar reads "SimulCAT [Vm~]". The menu bar includes "File" and "About". The main window has three tabs: "STEP1: Examinee / Item Data", "STEP2: Item Selection", and "STEP3". The "STEP1" tab is active, displaying the "Examinee Characteristics" dialog box. This dialog box contains the following fields and options:

- Number of Examinees:** A text input field containing the value "10000".
- Distribution:** A section with three radio buttons: "Normal" (selected), "Uniform", and "Beta".
- Mean:** A text input field containing the value "0".
- SD:** A text input field containing the value "1".
- Regenerate for each replication:** An unchecked checkbox.
- Generate True Scores (Theta):** A large green button at the bottom of the dialog box.

- If you have a special theta distribution, you could generate it using another program and import the thetas.



1...mhtropA...loo9 msl...noislumj2  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations



- Click generate true scores (Theta) and the theta values will be randomly generated and will appear in the Examinee( $\theta$ ) box.
- You can click the histogram button to view the distribution and various descriptive statistics.
- You also want to confirm the D value against your IRT program

[Vm~] SimulCAT

File About

STEP1: Examinee / Item Data STEP2: Item Selection STEP3: Test Administration

**Examinee Characteristics**

Number of Examinees

Distribution

Normal  Uniform  Beta

Mean  SD

Regenerate for each replication

**Generate True Scores (Theta)**

**Item Characteristics**

Number of Items

Distribution

Par.a  Mean  SD

Par.b  Mean  SD

Par.c  Mean  SD

Content ID

Scale to normal metric (scaling factor D=1.702)

Add to the previous item set

**Generate True Item Parameters**

**Examinee( $\theta$ )**

ID	$\theta$
1	0.081
2	-0.623
3	1.561
4	0.440
5	-3.252
6	0.578
7	-0.246
8	0.873
9	2.965
10	-0.739
11	0.407
12	-1.435
13	0.671
14	1.415
15	-0.514
16	1.185
17	0.238
18	-1.229
19	-1.193
20	-0.684
21	-0.173
22	-2.093

N = 10000  
Mean = -0.009  
SD = 1.006

Histogram



1...mltltropA...loo9 msl...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- Next, move to the STEP2: Item Selection tab.
- Here you will see options for many different types of algorithms.
- In the following slides, we will give a high-level overview of each.

SimulCAT [Vm~] SimulCAT

File About

STEP1: Examinee / Item Data | **STEP2: Item Selection** | STEP3: Test Administration

**Item Selection Criterion**

- Maximum Fisher Information (MFI)
- a- Stratification (Number of Strata:  )  b-blocking
- Matching b-Value
- Random Selection
- Fixed Sequence
- Interval Information Criterion (IIC)
- Likelihood Weighted Information (LWI)
- Kullbak-Leibler Information (KLI: Global Information) Constant 'c'
- Gradual Maximum Information Ratio (GMIR)
- Efficiency Balanced Information (EBI)

**Item Exposure Control**

- No Exposure Control
- Fade Away Method (FAM: Target Exposure Rate:  )
- Rangesque (Randomly select an item among  best items)
- Probabilistic Approach with Simulations
  - Sympson and Hetter Method (SHM)
  - Unconditional Multinomial Method (UMM)
  - Conditional Multinomial Method (CMM)
- Load Item Exposure Parameter Data
  - 
  - Derive Item Exposure Parameters after  iterations
    - Target Exposure Rate:
    - (only for CMM, number of Intervals:  between  and )
- Cumulative Item Usage Criterion for Retirement
- Item Usage Update (only for FAM and/or Retirement)
  - Real Time
  - After Each Test Time Slot
  - Daily Basis

**Test Length**

- Fixed Length  Items
- Variable Length
  - Terminate when SEE becomes smaller than
  - Terminates when the change in interim estimates becomes smaller than  for the last  items
  - Minimum  Items
  - Maximum  Items
  - Expected Length  Items (only for 'a-strat.' and 'GMIR')

**Content Balancing**

- None
- By Script
- By Weight

Ready | Elapsed Time 0:00:00:01





# Monte Carlo Simulations



- **Item Selection Methods**

- **Maximum Fisher Information (MFI):** Selects the item that would provide the maximum information value for the examinee's current theta estimate.
- **a-stratification:** The items are divided into different strata based on their a parameters. At the beginning of the CAT, the least discriminating stratum is used. The item with the closest b value to the examinee's current theta estimate is chosen. You can specify the number of strata. (Chang & Ying, 1999)
  - **b-blocking:** In addition to being formed based on a parameters, the strata are also designed so that they have balanced b parameter distributions. (Chang et al., 2001)



# Monte Carlo Simulations



- Matching-b: Selects the item with an item difficulty  $b$  parameter that is closest to the current theta estimate.
- Randomization: This is a control condition in which item are randomly selected.
- Interval Information Criterion: A variation of the MFI method (Veerkamp & Berger, 1997). This approaches sets up a confidence interval around the theta estimate and takes an averages of the information function across the interval.
- Likelihood Weight Information Criterion: Another variation of the MFI method (Veerkamp & Berger, 1997). This approach takes the likelihood function based on the previously administered items and uses it to weight the information function across the theta scale. The values are then summed.



# Monte Carlo Simulations



- Kullback-Leibler (1951) Information (Global Information): (Chang & Ying, 1996). Kullback-Leibler information is a type of information function. The global information method takes a moving average of that function.
- Gradual Maximum Information Ratio: Both MFI and the effective efficiency (i.e., how well the potential information for an item is realized for a given theta estimate) for items are computed. Both are summed together and the sum is used to select items. At the beginning of the CAT, more weight is given toward efficiency and at the end of the CAT more weight is given to MFI. (Han, 2009)



1. Multi-Item...  
2. 3-PLM...  
a, b, c... Monte Carlo

# Monte Carlo Simulations



- **Efficiency-Balanced Information:** Similar to the Gradual Maximum Information Ratio, except that item efficiency and MFI are evaluated across the theta estimate interval (which depends on standard errors; Han, 2010).
- **Item Exposure Methods**
- **Randomesque:** rather than selecting the single best item, the best  $k$  items are identified and one of those is randomly selected to administer to the examinee (Kingsbury & Zara, 1989).



# Monte Carlo Simulations



- **Sympson-Hetter (1985):** A target probability that an item is administered is established. A single random number between 0 and 1 is computed. The items are rank-ordered by the item selection method. Starting from the top, the target probability is compared to the random number. If the random number is smaller, the item is given, if not, the next item is considered.
- **Unconditional Multinomial Method:** A variation of Sympson-Hetter. A multinomial distribution is computed and then compared to the random number. (Stocking & Lewis, 1995)



# Monte Carlo Simulations



- **Conditional Multimonial Method:** A variation of the Unconditional Multimonial Method that applies the item exposure controls to different regions of the theta distribution. This approach would control item exposure for groups of examinees with similar ability. (Stocking & Lewis, 1995)
- **Fade-away method:** This approach requires continually tracking and updating observed item exposure data. The variable that item selection is based on (e.g., MFI) is weighted by the target exposure rate divided by the actual exposure rate. Items that are frequently used tend to fade-away from being administered.



...mlitrogA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Test Length**
- Fixed length: Each examinee receives the same fixed number of items.
- Variable length: Instead of administering the same number of items to each examinee, the actual number administered depends on the stability of the theta estimate.
  - Can stop testing when the standard error of the theta estimate is lower than a user-specified value
  - Can stop when changes in the theta estimates are below a user-specified value.
  - Can also specify a minimum and maximum number of items.



1...mltltrogA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Content-Balancing Methods**
- Script method: Several scripts are created representing different content areas.
- By weight/Constrained CAT method: An item is selected from a content area and each content area has a target percentage. The content area whose actual percentage of items administered and target percentage is the most different is selected. (Kingsbury & Zara, 1989)





1...mhtitopA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- At this point, determine which methods you want to explore.
- You can only test one method at a time.
- However, you can create syntax and cue files to run multiple methods.

SimulCAT [Vm~] SimulCAT

File About

STEP1: Examinee / Item Data | **STEP2: Item Selection** | STEP3: Test Administration

**Item Selection Criterion**

Maximum Fisher Information (MFI)

a- Stratification (Number of Strata: )  b-blocking

Matching b-Value

Random Selection

Fixed Sequence

Interval Information Criterion (IIC)

Likelihood Weighted Information (LWI)

Kullbak-Leibler Information (KLI: Global Information) Constant 'c'

Gradual Maximum Information Ratio (GMIR)

Efficiency Balanced Information (EBI)

**Item Exposure Control**

No Exposure Control

Fade Away Method (FAM: Target Exposure Rate: )

Rangesque (Randomly select an item among  best items)

Probabilistic Approach with Simulations

Sympson and Hetter Method (SHM)

Unconditional Multinomial Method (UMM)

Conditional Multinomial Method (CMM)

Load Item Exposure Parameter Data

Derive Item Exposure Parameters after  iterations

Target Exposure Rate:

(only for CMM, number of Intervals:  between  and )

Cumulative Item Usage Criterion for Retirement

Item Usage Update (only for FAM and/or Retirement)

Real Time  After Each Test Time Slot  Daily Basis

**Test Length**

Fixed Length  Items

Variable Length

Terminate when SEE becomes smaller than

Terminates when the change in interim estimates becomes smaller than  for the last  items

Minimum  Items

Maximum  Items

Expected Length  Items (only for 'a-strat.' and 'GMIR')

**Content Balancing**

None

By Script

By Weight

Ready | Elapsed Time 0:00:00:01



1...mltltropA...loo9 msl...noislumiz  
2, 3-PLM...θ, a, b, c...Monte Carlo

## Monte Carlo Simulations



- Next move to Step3.
- Here you can indicate test administration and scoring conditions.
- We discuss each on the following slides.

The screenshot shows the SimulCAT software interface at Step 3: Test Administration. The window title is "[Vm~] SimulCAT". The interface is divided into several sections:

- Score Estimation:** Includes radio buttons for Maximum Likelihood Estimation (MLE), MLE with Fences (MLEF), Bayesian Maximum a Posteriori (MAP), and Bayes Expected a Posteriori (EAP). It also has input fields for Lower Fence (-3.5), Upper Fence (3.5), Prior Mean (0), and Prior SD (1). There is an "Initial Score Value" section with options for Fixed Value, Randomly Chosen Value Between (-0.5 and 0.5), and From a Data File.
- Options:** Includes checkboxes for "Limit (Truncate) Range of Estimates" (Lower Bound -3, Upper Bound 3), "Limit the Estimate Jump by" for the first items, and "Use MLE for the Final Estimate".
- Test Administration:** Includes input fields for "Number of Examinees for Each Test Time Slot" and "Number of Test Time Slots per Day".
- Pretest Item Administration:** Includes a radio button for "None" and an "Administer" checkbox with an "Open Pretest Item File" button.
- Extras:** Includes checkboxes for "Generate Replication Data Sets", "Previous Item Usage Data", "Fixed Seed Value", and "Item Pool with DIF/Drift". It also has a "Run Simulation" button.
- Outputs:** Includes checkboxes for "Save Response Strings and Item IDs", "Save Item Use (Exposure) (in a \*.SCU file)", "Save Interim SEEs and TIFs", "Save Newton-Raphson Iteration Info", "Save Interim Theta Estimates", "Save Full Response Matrix", and "Save True Interim SEEs and TIFs".
- Log/Message:** A large text area for logging simulation results.



# Monte Carlo Simulations



- **Score Estimation**
- These are essentially the classic theta estimation approaches: Maximum likelihood estimation (MLE), Bayesian maximum a posteriori (MAP), and Bayes expected a posteriori (EAP).
  - MLE with Fences is MLE but artificial items with fixed responses are added to allow estimating theta for abnormal response patterns and those that consist of all 0s or 1s. The fence values are user-specified b parameters (see Han, 2016, for more information).
- We recommend using whatever approach you used to estimate the item parameters



1...multitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Initial Score Value**
- Fixed value: All examinees start the CAT with the same initial theta (often set to 0).
- Randomly chosen value between a user-specified interval
- From a data file: In this option, you can supply starting values for each examinee. This might be used if you had a previous test, pre-test, or other information for choosing a starting value.



1...mltrpA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Options**
- You can also truncate the range of estimates to a specific interval (e.g., if a theta estimate was -4, you could truncate it to -3), limit how much the theta estimates change and apply that limit to the first k items, or use MLE as the final theta estimate.
- **Test Administration**
  - Some of the item exposure controls depend on item exposure rates that are computed daily; here you can customize how many examinees are in a time slot and how many time slots are in a day.



...mlitroqA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Pretest item administration**
- You can simulate administering a pretest to examinees.
- **Extras**
  - It's also possible to conduct multiple replications, study differential item functioning (DIF) and item drift, consider previous item usage, and fix a random number seed value.
    - We recommend fixing the random number seed value for each run so that you could replicate the exact results if you later needed to.



1...mlitrogA...loo9 msl...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



## ■ Outputs

- You can save all types of output for each run. We found it useful to save the item use exposure (.SCU) file (so that you can track item exposure without having to do too much data manipulation).
- At a bare minimum you also need the true and final theta estimates.
- You could, of course, check all the boxes and save everything. However, if you are testing many conditions and have a large sample size, file sizes and disk space can quickly become an issue.



1...mltltropA...loo9 mell...noitlumiz  
2, 3-PLM...θ, a, b, c...Monte Carlo

## Monte Carlo Simulations



- Finally, click run simulation.

The screenshot shows the SimulCAT software interface with the following sections:

- Score Estimation:** Includes radio buttons for Maximum Likelihood Estimation (MLE), MLE with Fences (MLEF), Bayesian Maximum a Posteriori (MAP), and Bayes Expected a Posteriori (EAP). It also has input fields for Lower Fence (-3.5), Upper Fence (3.5), Prior Mean (0), and Prior SD (1). There is an 'Initial Score Value' section with options for Fixed Value, Randomly Chosen Value Between (-0.5 and 0.5), and From a Data File.
- Options:** Includes checkboxes for 'Limit (Truncate) Range of Estimates' (Lower Bound -3, Upper Bound 3), 'Limit the Estimate Jump by' (for the first items), and 'Use MLE for the Final Estimate'.
- Test Administration:** Includes input fields for 'Number of Examinees for Each Test Time Slot' and 'Number of Test Time Slots per Day'.
- Pretest Item Administration:** Includes a radio button for 'None' and an 'Administer' checkbox with an 'Open Pretest Item File' button.
- Extras:** Includes checkboxes for 'Generate Replication Data Sets', 'Previous Item Usage Data', 'Fixed Seed Value', and 'Item Pool with DIF/Drift'. It also has an 'Item Review Options' section with radio buttons for 'None', 'Item Pocket', 'Modified EoT', 'Moving Window', and 'End-of-Test (EoT)'. There is an 'Outputs' section with checkboxes for 'Save Response Strings and Item IDs', 'Save Item Use (Exposure) (in a \*.SCU file)', 'Save Interim SEEs and TIFs', 'Save Newton-Raphson Iteration Info', 'Save Interim Theta Estimates', 'Save Full Response Matrix', and 'Save True Interim SEEs and TIFs'.
- Output File:** Includes an 'Output File' input field and a 'Browse' button.
- Log/Message:** A large text area for logging and messages.

The 'Run Simulation' button is highlighted with a red oval.





1...multitropA...loo9 mell...noislumig  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Designing your Monte Carlo Simulation**
- You will have to think strategically about which conditions to study. The design can get out of hand pretty quickly.
  - For one CAT, we thought about testing almost all of the options but then realized it would have taken 10 years to run the simulation.
  - We recommend thinking carefully about each option and choosing those you are comfortable with considering psychometrics, programming feasibility, and face validity.
    - For example, starting with a randomly-chosen initial theta value might help with item exposure, but some examinees might perceive this as unfair and arbitrary.



1...multitopA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Analyzing the Results from your Monte Carlo Simulation**
- We recommend computing a number of summary statistics for each condition and then comparing the results.
- **Mean Bias:** This is the final theta estimate minus the true theta estimate averaged across all examinees. It will tell you if the algorithm, on average, tends to over or underestimate theta.
- **Mean |Bias|:** It's also helpful to convert the mean bias values to positive values to allow you to sort the results for each condition and find those that are closest to zero.



1...multitropA...loo9 mell...noislumig  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Analyzing the Results from your Monte Carlo Simulation**
- Mean Square Error:  $(\text{theta final} - \text{true theta})^2$ .
- Root Mean Square Error: Square root of MSE
- Correlation: correlation between final theta estimates and true thetas.
- Reliability: The correlation squared.





# Monte Carlo Simulations

## ■ Analyzing the Results from your Monte Carlo Simulation

- Chen et al. (2003) average test overlap rate:  $\bar{T} = \frac{S^2 + \mu^2}{\mu}$ 
  - Where  $S^2$  is the variance of the item exposure rates (which are the proportion of examinees that received an item),  $\mu$  is the mean of the item exposure rates.
  - A value of 40% means that, on average, two examinees have 40% of their items in common.
- Number of underused items (i.e., number of items administered to less than .05 [5%] of examinees)
- Percentage of items that are underused



1...multitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Monte Carlo Simulations



- **Analyzing the Results from your Monte Carlo Simulation**
- Scaled  $\chi^2$  (Chang & Ying, 1999): This compares the observed item exposure to an ideal/expected uniform item exposure rate. The ideal/expected rate is the length of the CAT divided by the total number of items in the item pool.
- Essentially, you will compare these statistics across different conditions and decide on a final algorithm.
  - You might also make refinements and run additional conditions.



1...mlitropA...loo9 mell...noislumiz  
2, 3-PLM... $\theta$ , a, b, c...Monte Carlo

# Scaling/Equating/Metric Issues



- The CAT will generate a final  $\theta$  estimate for each examinee.
- You could use that the final  $\theta$  estimate as the operational score.
  - It is on a z-scale and has positive, negative, and zero values. This makes it more difficult to explain to non-psychometricians.
  - It can be useful to place the  $\theta$  estimates onto the same scale as a static test.
  - Oftentimes, you are replacing a static test with a CAT.
  - The static test's metric/scale might have been used in validation and standard-setting studies.



1...multitropA...loo9 mell...noislumj2  
2, 3-PLM... $\theta$ , a, b, c...Monte Carlo

# Scaling/Equating/Metric Issues



- Hambleton, Swaminathan, and Rogers (1991, p. 84-7) explain how to place  $\theta$  on a number correct scale for a static test.
  - You first need to obtain the item parameters for the static test.
  - Next, determine the desired level of precision in the equating (e.g., do you want to equate  $\theta$  in increments of 0.1, 0.01, etc.).



# Scaling/Equating/Metric Issues



- Create a spreadsheet with different  $\theta$  values in each row and  $P$  (probability of correct response) for each item in columns. Use the 1-, 2-, or 3-PLM equation to compute the  $P$ s. Next sum the  $P$ s across the items to give the number of correct items ( $T$ , using Hambleton et al.'s notation). You can divide this by the total number of items to give the proportion correct ( $\pi$ ), and convert to a percentage if desired.
- On the next slide, we give an example for a fictitious 10-item test with  $\theta$  in unit increments.





Monte Carlo  
2, 3-PLM...θ, a, b, c...

## Scaling/Equating/Metric Issues



Theta transformation for CAT Tutorial.xlsx - Excel

File Home Insert Page Layout Formulas Data Review View ACROBAT Tell me what you want to do... CUCINA, JEFFREY Share

B2  $=B\$12+((1-B\$12)*((2.718^{(1.7*B\$10*(\$A2-B\$11))})/(1+(2.718^{(1.7*B\$10*(\$A2-B\$11))}))))$

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Θ	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	T	π
2	-3	0.214	0.231	0.102	0.259	0.209	0.164	0.325	0.173	0.246	0.113	2.036	20%
3	-2	0.369	0.235	0.118	0.269	0.288	0.211	0.335	0.234	0.253	0.131	2.445	24%
4	-1	0.677	0.263	0.193	0.311	0.528	0.403	0.389	0.392	0.315	0.211	3.683	37%
5	0	0.903	0.431	0.454	0.450	0.826	0.750	0.577	0.646	0.614	0.463	6.113	61%
6	1	0.978	0.800	0.805	0.710	0.958	0.942	0.842	0.856	0.923	0.795	8.610	86%
7	2	0.995	0.968	0.957	0.906	0.991	0.989	0.964	0.953	0.991	0.951	9.666	97%
8	3	0.999	0.996	0.992	0.976	0.998	0.998	0.993	0.986	0.999	0.990	9.928	99%
9	Item Parameters												
10	a	0.923	1.234	1.01	0.875	0.955	1.022	0.998	0.735	1.31	0.958		
11	b	-1.3	0.5	0.25	0.7	-0.8	-0.5	0.3	-0.279	0.021	0.257		
12	c	0.159	0.231	0.099	0.256	0.187	0.153	0.322	0.145	0.245	0.109		
13													



Monte Carlo  
3-PLM...θ, a, b, c...

# Scaling/Equating/Metric Issues



Theta transformation for CAT Tutorial.xlsx - Excel

File Home Insert Page Layout Formulas Data Review View ACROBAT Tell me what you want to do... CUCINA, JEFFREY Share

B2 =B\$12+((1-B\$12)\*((2.718^(1.7\*B\$10\*(\$A2-B\$11)))/(1+(2.718^(1.7\*B\$10\*(\$A2-B\$11))))))

	A	B	C	D	E	F	G	H	I	J	K	L	M
1	θ	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	T	π
2	-3	0.214	0.231	0.102				0.325	0.173	0.246	0.113	2.036	20%
3	-2	0.369		0.118				0.335	0.234	0.253	0.131	2.445	24%
4	-1	0.677	0.265					0.389	0.392	0.315		3.683	37%
5	0	0.903	0.431					0.577	0.646	0.614	0.463	6.113	61%
6	1	0.978	0.800	0.805				0.815	0.855	0.855	0.755	8.610	86%
7	2	0.995	0.968	0.957				0.915	0.955	0.955	0.955	9.66	97%
8	3	0.999	0.996	0.992				0.985	0.995	0.995	0.995	9.9	99%
9	Item Parameters												
10	a	0.923	1.234	1.01				0.923	1.234	1.01			
11	b	-1.3	0.5	0.25				-1.3	0.5	0.25			
12	c	0.159	0.231	0.099				0.159	0.231	0.099			
13													

The values in cells B2:K8 are the Probability of correct responses using the 3-PLM equation, which is typed in the bar above.

The values in cells L2:L8 are the Sum of the values in columns B through K

This is the number correct converted to a percentage.



# Scaling/Equating/Metric Issues



- A few notes:
  - The number correct scores and  $\theta$ s are monotonically related (not necessarily linear).
  - You can apply this procedure to compute number correct scores for item that the examine didn't receive, provided that you have the  $\theta$ s and item parameters (and both are on the same metric).
  - This is actually a transformation of true  $\theta$  scores to true number correct scores.
  - This procedure is a derived from the summation of item characteristic curves to yield a test characteristic curve. However, the curves are usually for the static test items, not the entire CAT item pool; therefore, the curves we discussed earlier and that are produced in SimulCAT (for example) are not relevant here.



1...multitrogA...loo9 mell...noitlumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Experimental Item Collection



- After CAT is implemented, you may desire to continue collecting data on experimental items
  - This will help you to develop additional CATs
  - Might need to add new items if content domain changes over time (e.g., some aspects of job knowledge might change over time).
  - You might also need to recalibrate existing items to control for item drift.
    - Suppose you developed 5 CATs each with their own pool. At some point in the future, you might consider stopping use of one of those CATs and readministering the items as experimental items in order to recalibrate the item parameters.



...mlitrogA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Experimental Item Collection



- There are some strategies for collecting experimental items
  - Administered a static test in addition to the CAT
    - Need to consider order effects and fatigue
    - Examinees likely need to be told that one of the tests contains unscored experimental items.
      - They might be able to identify which test is the non-CAT static experimental one and then reduce their motivation when completing items.



...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Experimental Item Collection



- There are some strategies for collecting experimental items
  - Embed experimental items within the CAT
    - Experimental items could be in fixed or random locations
    - We suggest avoiding placing experimental items at the very beginning or end of the CAT
      - Examinees often know that the first item should be of average ability and the last should be near their own ability levels. They might panic if the first item is very hard or if the last item is very easy.
      - Examinees also might be gauging their performance throughout the CAT and could possibly identify or be thrown off by experimental items that are far from interim theta.



Monte Carlo  
a, b, c...  
2, 3-PLM...  
1...multitropA...loo9 mell...noislumiz

# Experimental Item Collection



- There are some strategies for collecting experimental items
  - Collect experimental items outside of CAT environment (e.g., using a static test that is part of the testing program).
- In most cases, you will have a set of unscored static test items, CAT item responses, and thetas.
- In the literature, the preferred method seems to be to have enough usable (i.e., those that will survive item analysis) static test items to conduct a calibration and obtain item parameters and thetas. (This can range from 10 to 60 items depending on which “rules of thumb” you choose.)



# Experimental Item Collection



- The static items are calibrated separately from the CAT responses and could be on a different scale.
- Pommerich & Segall (2003) developed a transformation which works well according to simulation studies for placing static test (Test 2) parameters on same scale as CAT (Test 1); note that only a and b, but not c, need transformation:

$$a_{\text{Transformed, Test 2}} = a_{\text{Test 2}}/A$$

$$b_{\text{Transformed, Test 2}} = A(b_{\text{Test 2}})+B$$

Where:

$$A = \sigma_{\text{Test1}}/\sigma_{\text{Test2}}$$

$$B = \mu_{\text{Test1}} - A(\mu_{\text{Test2}})$$





# Experimental Item Collection



- Other approaches include
  - Attempting to calibrate CAT responses with static items.
    - This yields a sparse data matrix with a lot of missing data in the CAT responses.
  - Using BILOG-MGs external  $\theta$  command (which allows you to place new item parameters on same scale as an existing set of  $\theta$ s)
    - However, we could not locate any technical details on how this is done in BILOG-MG or how well it works



1...mlitrogA...loo9 mell...noislumiz  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Experimental Item Collection



- Other approaches include
  - Treating CAT  $\theta$  estimates as known and estimating item parameters when  $\theta$  is known.
    - However,  $\theta$ s are actually estimates with error and not known and most IRT software packages don't implement this option.
  - Including a set of static anchor items with known item parameters and equating the experimental items to the anchor items using IRTEQ (for example).





1...mlitrogA...loo9 mell...noitlumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# Creating Instructions for Developers



- Recommended information to include in the instructions (continued):
  - Number of quadrature points (e.g., 20, 40)
  - Instructions for conducting random selections of items
  - Equating instructions to other test forms
  - Location of experimental items (for pretesting other items)
  - Output file formats



# Creating Instructions for Developers



- A critical step in developing the instructions is to think strategically about which variables should be computed and recorded when the CAT is administered.
- This is valuable in case the CAT is being used in a high-stakes situation where scores may be challenged.
- It is recommended to audit the programmed CAT algorithm to ensure no mistakes were made in the programming and the ability to do this depends largely on the variables that are recorded and outputted.



1...mltltrogA...loo9 mell...noislumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# CAT User Testing



- After a CAT algorithm is chosen and implemented, it is recommended that careful user testing is conducted to ensure that the CAT algorithm is functioning properly.
- User testing ensures no errors were made during the programming process. Ensuring the CAT algorithm is error-free provides protection for organizations and appropriate legal defensibility.
- We recommend creating a very large (e.g., 500,000) number of test cases and running them through the CAT algorithm using a Monte Carlo simulation to ensure that all possible item iterations are verified.



1...multitropA...loo9 mell...noitlumj2  
2, 3-PLM...0, a, b, c...Monte Carlo

# CAT User Testing



- Several aspects of CAT output that we have identified as crucial to user testing include:
  - Interim theta estimates (i.e., the ability estimates between item administrations).
  - Item selection correctness (i.e., whether the correct items are being selected based on the interim thetas).
  - Final thetas (i.e., the final ability estimates of the test).
- Interim theta estimates and item selection correctness are important because theoretically a test taker could have a correct final ability estimate, but an incorrect item iteration or interim theta estimates.



...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# CAT User Testing



- Final thetas should be examined for correctness to the last decimal place and rounding errors, as these issues may be problematic around where a test cut score is set.
- Additionally, there may be legal or other challenges that require verifying or proving that the test was administered correctly, and these aspects of the test may need to be verified or produced.





...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Implementing CATs



- When implementing a CAT, it is important to update any applicable examinee communications
  - Provide a description of the new CAT
  - Give insight on what examinees can expect and how they should prepare for taking a CAT
  - Update relevant website(s)
    - Provide updated information on CAT
    - Indicate when the CAT will be implemented
- The CAT developer may also need to work with stakeholders (e.g., HR staff, recruiters) so that they are aware of the new test and implementation guidelines



...mlitropA...loo9 mell...noitlumiz  
2, 3-PLM...θ, a, b, c...Monte Carlo

# Implementing CATs



- CATs can be either developed from scratch or from an existing static test
  - If a psychometrician is converting a static test to a CAT, then he or she will need to decide if it is necessary to discontinue use of existing static test scores upon implementation of CAT
    - A decision will need to be made whether applicants or candidates with passing scores on the static test need to take the adaptive version or if they will be allowed to keep their passing score.
  - This will not be a concern if the CAT is being developed as a brand new test



...mlitropA...loo9 mell...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Implementing CATs



- After the CAT is implemented, it is important to confirm the algorithm was implemented correctly
  - Checking live data is the only way to know if the CAT is functioning as intended
    - Is the CAT selecting the correct items for test takers?
    - Are incorrect and correct response options being scored correctly?
    - Are the stopping rules working correctly?
    - Are calculations for interim and final thetas correct?
    - Do final thetas match any equated scores?
    - Is the scoring for the overall test correct?



1...mltirogA...loo9 msl...noislumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Implementing CATs



- To minimize any potential issues with the algorithm that could arise after the implementation of the CAT, it is recommended that a plan be created for checking the data prior to implementation
  - Indicate what part(s) of the algorithm and scoring you need to check
  - Provide the steps for how those checks will be carried out



Monte Carlo  
a, b, c...  
2, 3-PLM...  
1...multitropA...loo9 mell...noislumiz

# Lessons Learned/Things to Consider



- CAT provides a number of important advantages over traditional testing
- CAT can provide more accuracy because each examinee is given a unique test that is tailored to his or her ability level
  - Questions that do not provide sufficient information about the examinee's proficiency are avoided
    - For example, really easy questions are typically not asked to individuals with a very high ability level and vice versa
  - This can provide a higher level of precision across a wider range of ability levels



Monte Carlo  
a, b, c...  
2, 3-PLM...  
1...multitropA...loo9 mell...noislumj2

# Lessons Learned/Things to Consider



- In comparison to static tests, CATs can reduce testing length by more than 50%, while maintaining a comparable level of reliability
  - Fewer items are needed to achieve acceptable accuracy
- CAT testing has also been shown to improve test security
  - Since each test is unique to the examinee, this makes it more difficult for an individual to capture the entire pool of items
    - This drastically reduces the likelihood of widespread cheating and the subsequent need to redo an entire test, which would cost a significant amount of time and money



# Lessons Learned/Things to Consider



- There are some notable limitations to using the CAT method
- Unlike traditional testing, CAT is difficult to develop in that it requires the expertise of psychometricians to calibrate items using an algorithm based on an IRT model that accurately measures the ability level of examinees
  - Item calibration requires that extensive data be collected on a large item pool
    - Developing a sufficiently large item pool can take a lot of time and effort and can end up using a lot more resources than traditional testing



# Lessons Learned/Things to Consider



- Some subjects and skills cannot be accurately measured with CAT because IRT cannot be readily applied to those areas
- Using IRT models can also cause item constraints that result in an overly narrow selection of questions being presented to examinees
  - This is likely to occur if the content isn't balanced across different item difficulty levels.
  - These constraints can result in examinees completing sets of items that are broadly the same, thus losing the advantage over traditional tests
    - This presents a potential issue when CAT is used for knowledge and content-focused tests there can be blatant inequities when comparing the scores of examinees
    - Content balancing methods are a possible solution





# Lessons Learned/Things to Consider



- Some of the limitations of CAT can be addressed using another method of testing called linear-on-the-fly testing (LOFT)
- LOFT takes items from a very large item pool and constructs a unique exam for each examinee
  - This is done through a program that pseudo-randomly selects items so that examinees receive tests that are equivalent with respect to content and statistical characteristics, but with different items
    - Like CAT, this selection method typically uses IRT, but results in longer testing times and less precision than CAT



# Lessons Learned/Things to Consider

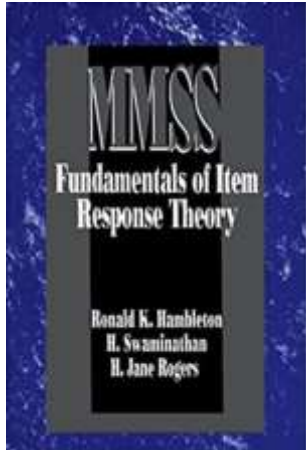


- Ultimately, when it comes to different modes of testing (e.g., static vs. CAT vs. LOFT), there will always be some trade-offs between the different methods
  - No single method will address every limitation
- When deciding whether to use methods such as CAT, it is important for test developers to weigh the advantages and disadvantages of each method and decide which approach works better for their situation



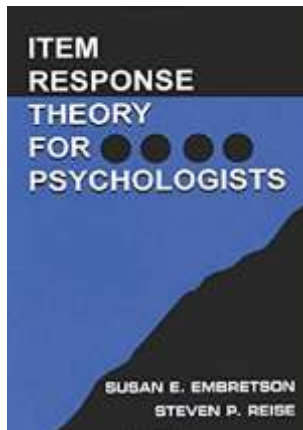
1...mltirogA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Resources/Recommended Reading



THE  
THEORY AND  
PRACTICE OF  
ITEM RESPONSE  
THEORY

R.J. de Ayala



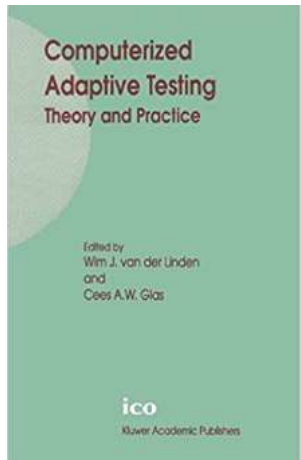
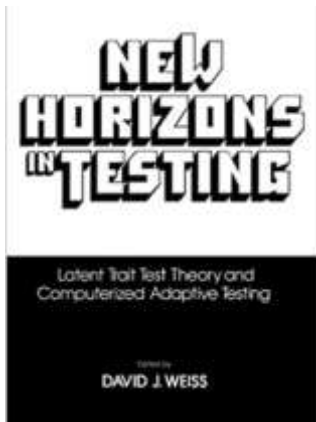
- Web-based item characteristic curve plotting
  - <https://demonstrations.wolfram.com/ItemCharacteristicCurves/>
- IRTEQ, SimulCAT, WinGen, and other free programs: <http://www.hantest.net/>
- Hambleton, R.K., Swaminathan, H, & Rogers, H.J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: SAGE Publications.
- de Ayala, R.J. (2009). *The theory and practice of item response theory*. New York: The Guilford Press.
- Hambleton, R. K., & Swaminathan, H. (2013). *Item response theory: Principles and applications*. New York: Springer Science & Business Media.
- Embretson, S. E., & Reise, S. P. (2013). *Item response theory for psychologists*. Mahwah, NJ: Lawrence Erlbaum Associates, Inc.



U.S. Customs and Border Protection

1...mlitropA...loo9 mell...noitlumiz  
2, 3-PLM...0, a, b, c...Monte Carlo

# Resources/Recommended Reading



- Weiss, D.J. (1983). *New Horizons in Testing: Latent trait test theory and computer adaptive testing*. New York, NY: Academic Press.
- Van der Linden, W.J., & Glas, C.A.W. (2000). *Computer adaptive testing: Theory and practice*. Dordrecht, The Netherlands: Kluwer Academic Publishers, Inc.
- Wainer, H. (2014). *Computer adaptive testing: A primer*. (2<sup>nd</sup> ed.). New York: Routledge.

