

# Item Response Theory and Computer-based Testing

1 March 2018

Seminar/Proseminar: Intelligente Lehr-/Lernsysteme  
Rainier Robles, WS 17/18

Name: Rainier Raymond Robles

Matrikelnummer: 5002400

- (a) Studienordnung: Bachelor Informatik (150 LP, Studienordnung 2014)
- (b) Modul: Wissenschaftliches Arbeiten in der Informatik 0086cA6.1
- (c) Modulprüfung: Vortrag (ca. 30 Minuten) mit anschließender Diskussion (ca. 10 Minuten)

Formen aktiver Teilnahme: Schriftliche Ausarbeitung, Teilnahme an den Diskussionen zum Vortrag

- I. Motivation
- II. Computer Adaptive Tests
- III. Item Response Theory
- IV. Case Study: SIETTE
- V. Issues and Challenges
- VI. Conclusions
- VII. References

# Motivation

"...to go back to a time when teachers could afford to assess each student orally by asking him/her a few well-selected questions..." (Conejo et al. 2004)

# Computer Adaptive Tests

"The basic notion of an adaptive test is to mimic automatically what a wise examiner would do." (Wainer and Mislevy, qtd. in Conejo et al.)

**Definition:** CATs are tests administered by a computer, where the selection of the items as well as the decision to terminate the test are based on the examinee's responses.

- Significantly shorter tests, leading to significantly shorter testing times
- More accurate estimation of test takers' proficiency levels
- Test takers' motivation improved
- Large item pools can be stored
- Multimedia content is supported



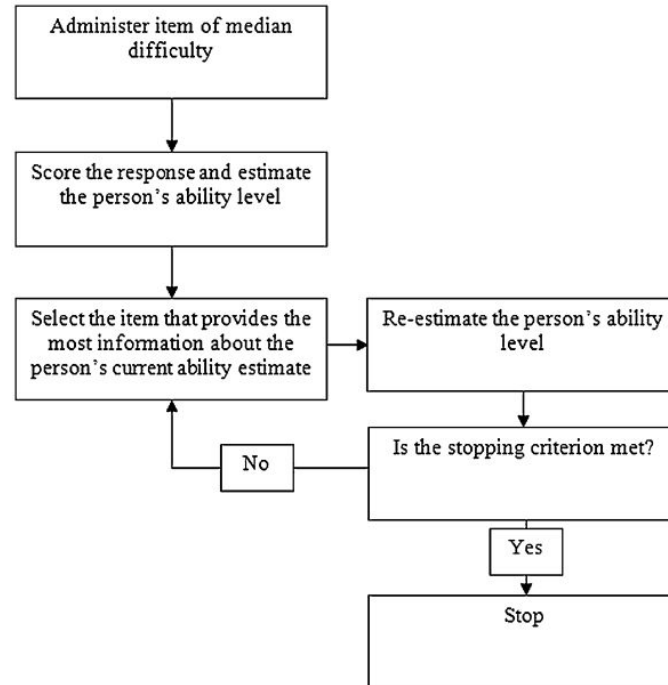


Fig. 1. CAT algorithm flowchart. (Wauters et al. 2010)

# Basic Elements in the Development of a CAT

- Response model
- Item pool
- Input proficiency level
- Item selection method
- Termination criterion

# Item Response Theory

## Classical Test Theory

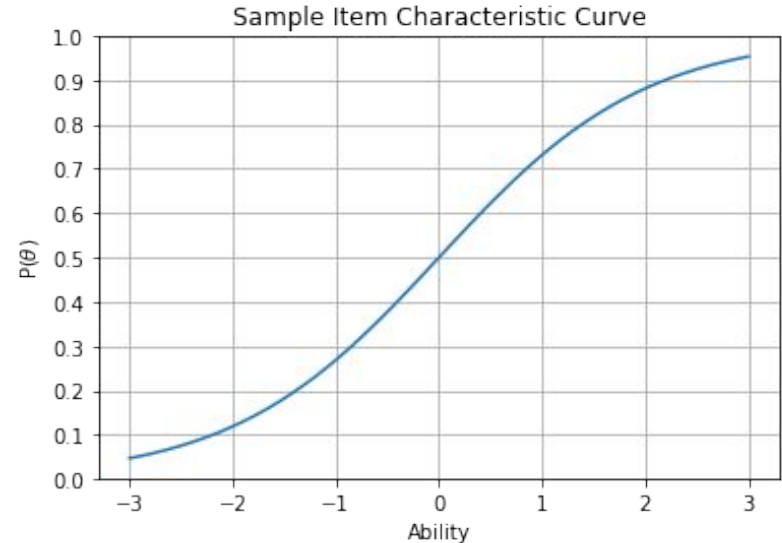
- Focus on aggregation of test responses as a total test score
- Generally produces a standardized achievement level
- Greater possibility to use free-response items in a test

## Item Response Theory

- Focus on individual responses to each test item
- Produces a probability of a test taker's proficiency level
- Mostly composed of multiple-choice items

# Item Characteristic Curve (ICC)

- "(T)he basic building block of item response theory." (Baker 2001)
- Plots probability that a test taker on a particular ability level can answer the question correctly
- Each item has its own ICC

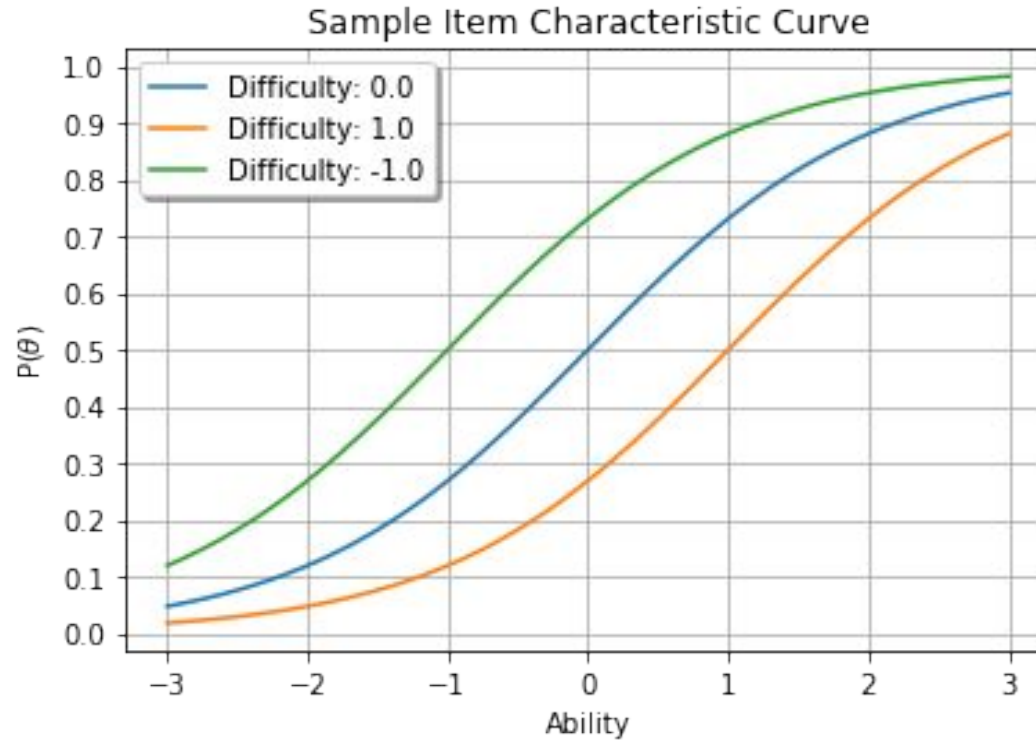


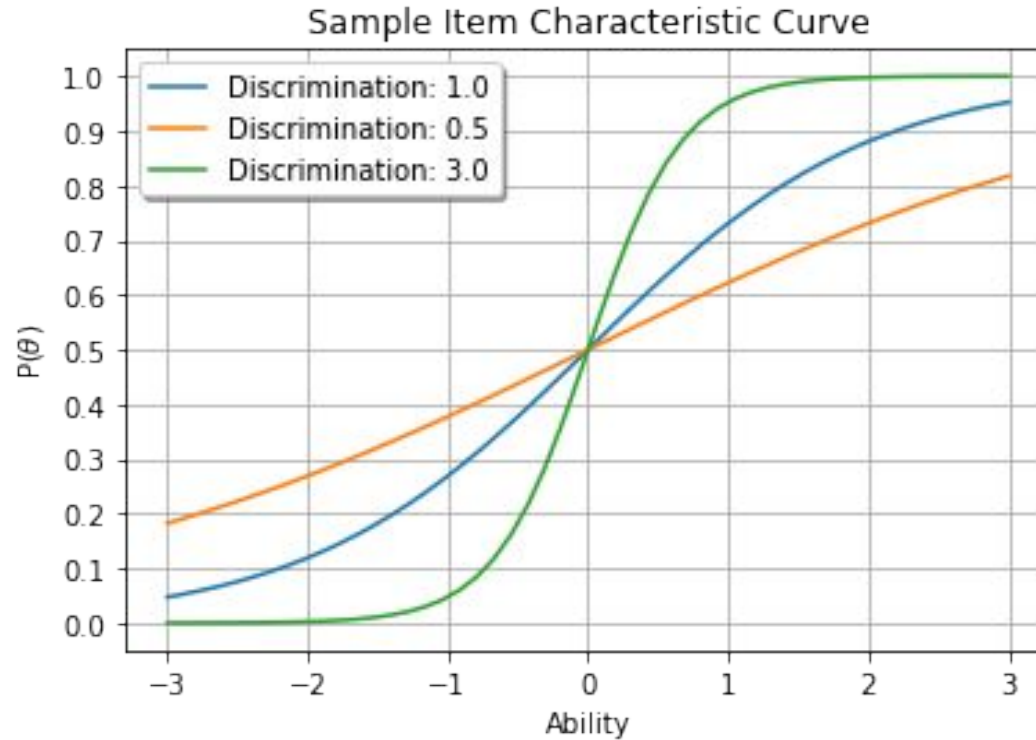
## 1. Item difficulty

Describes where an item "functions"; a location index.

## 2. Discrimination

Describes how well an item discriminates among test takers with abilities above and below the item's difficulty.







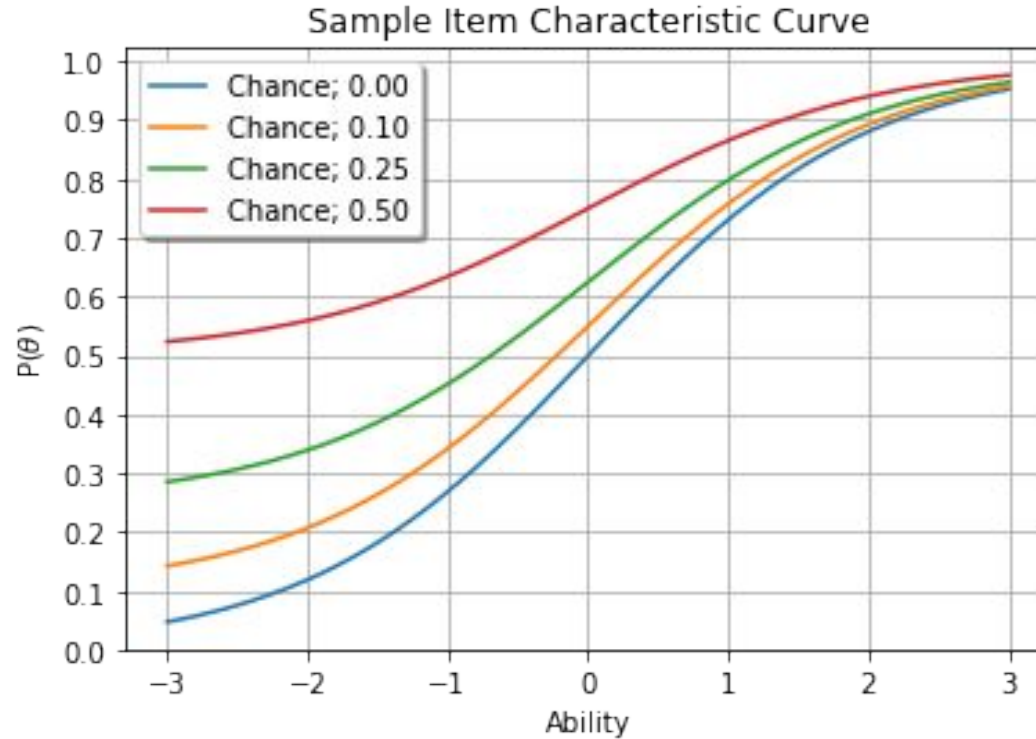
1. Two-Parameter Logistic Model
2. Rasch / One-Parameter Logistic Model
3. Three-Parameter (Logistic) Model

$$P(\theta) = c + (1 - c) \frac{1}{1 + e^{-a(\theta - b)}}$$

Where:

- $b$  = difficulty
- $a$  = discrimination
- $c$  = guessing parameter
- $\theta$  = ability level

# Three-Parameter Logistic Model



# Case Study: SIETTE

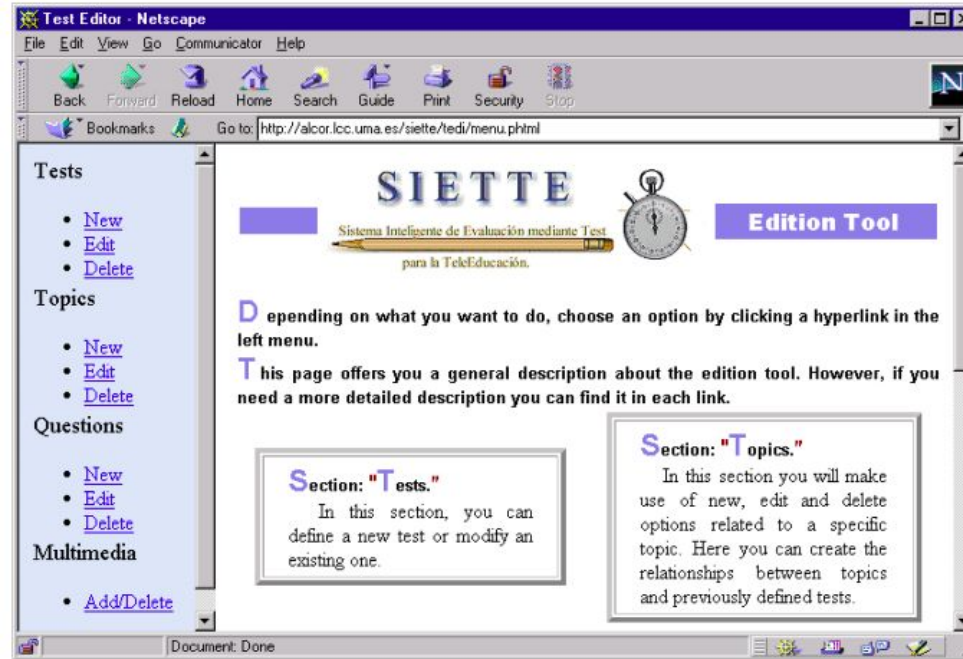


Fig. 2. SIETTE test editor homepage. (Conejo, et al. 2004)

- **S**ysteme de **E**valuación de **T**ests para la **T**ele**E**ducación (*Intelligent Evaluation System using Tests for TeleEducation*)
- Initially developed in 1998 by Antonia Rios as part of a master thesis
  - Department of Computer Science, University of Malaga, Spain
- Started as an Intelligent Tutoring System (ITS) module as part of the TREE (**T**Raining of **E**uropean **E**nvironmental Trainers and Technicians) project

# A Brief Overview of SIETTE

- Designed to be a reusable, domain-independent module
- Teachers can design and develop the parameters of the tests and items
- Students can take automatically generated adaptive tests
- Uses discrete ability levels

STUDENTS					
StudentID	TestID	Date	Level of Proficiency	Lower Confidence Level	Upper Confidence Level
John Smith	TREE	04/08/98	Level 1	0.9	1.1

KNOWLEDGE DISTRIBUTION					
StudentID	TestID	Level 0	Level 1	...	Level 10
John Smith	TREE	0.001	0.9	...	0.001

TOPIC DISTRIBUTION			
StudentID	TestID	TopicID	% Questions
John Smith	TREE	PINUS	40%
John Smith	TREE	ABIES	40%
John Smith	TREE	CEDRUS	20%

QUESTIONS POSED		
StudentID	QuestionID	AnswerID
John Smith	Q <sub>1</sub>	A <sub>1,1</sub>
John Smith	Q <sub>3</sub>	A <sub>3,2</sub>
John Smith	Q <sub>5</sub>	A <sub>5,1</sub>
John Smith	Q <sub>6</sub>	A <sub>6,3</sub>
John Smith	Q <sub>8</sub>	A <sub>8,4</sub>
John Smith	Q <sub>10</sub>	A <sub>10,5</sub>
John Smith	Q <sub>12</sub>	A <sub>12,3</sub>
John Smith	Q <sub>14</sub>	A <sub>14,1</sub>
John Smith	Q <sub>17</sub>	A <sub>17,2</sub>
John Smith	Q <sub>20</sub>	A <sub>20,1</sub>

Fig. 3. Example of a temporary student model in SIETTE. (Conejo, et al. 2004)



# Student State Before Answering

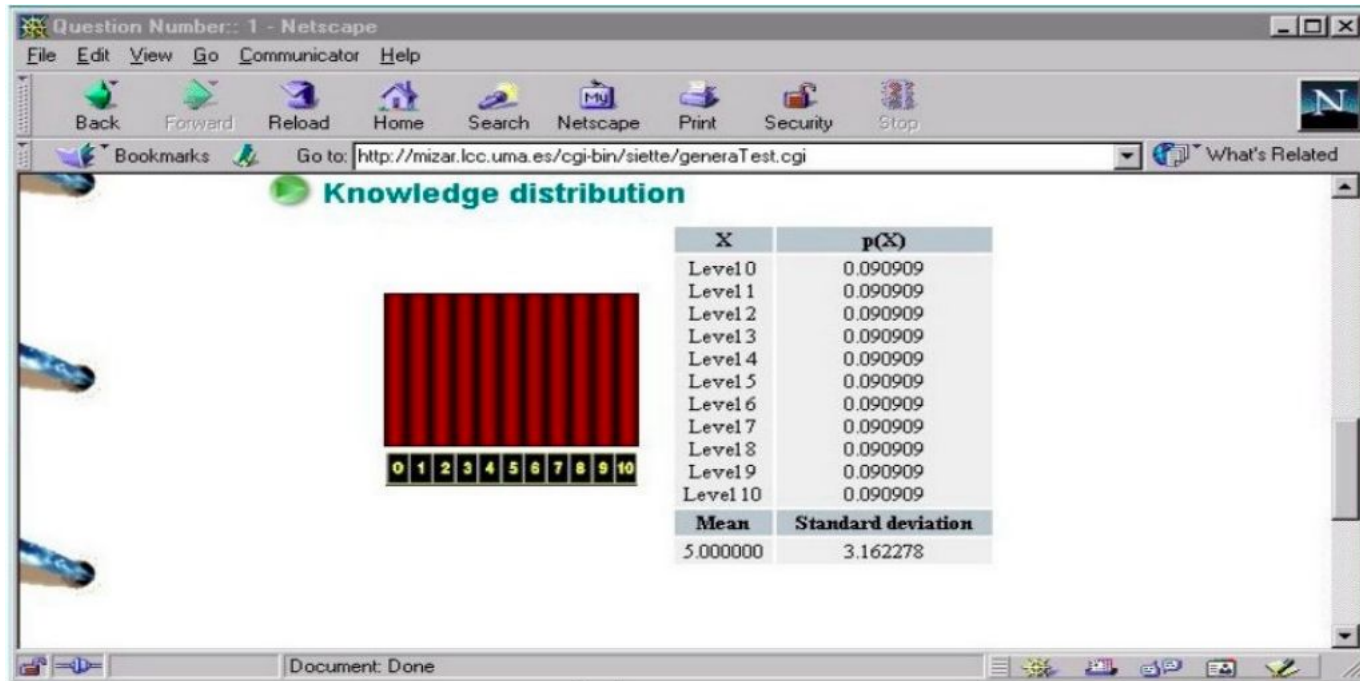


Fig. 4. Page showing estimated student levels before the first question. (Conejo, et al. 2004)

# Item Characteristic Curve in SIETTE

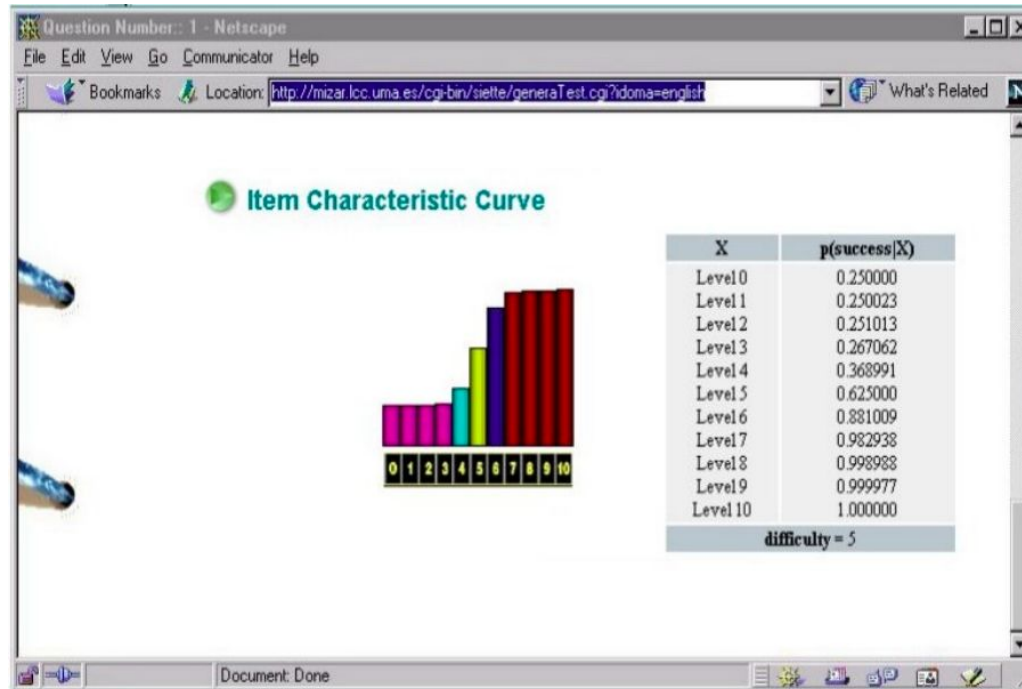


Fig. 5. ICC of an item in SIETTE. (Conejo, et al. 2004)

# Knowledge Distribution at the End of a Test

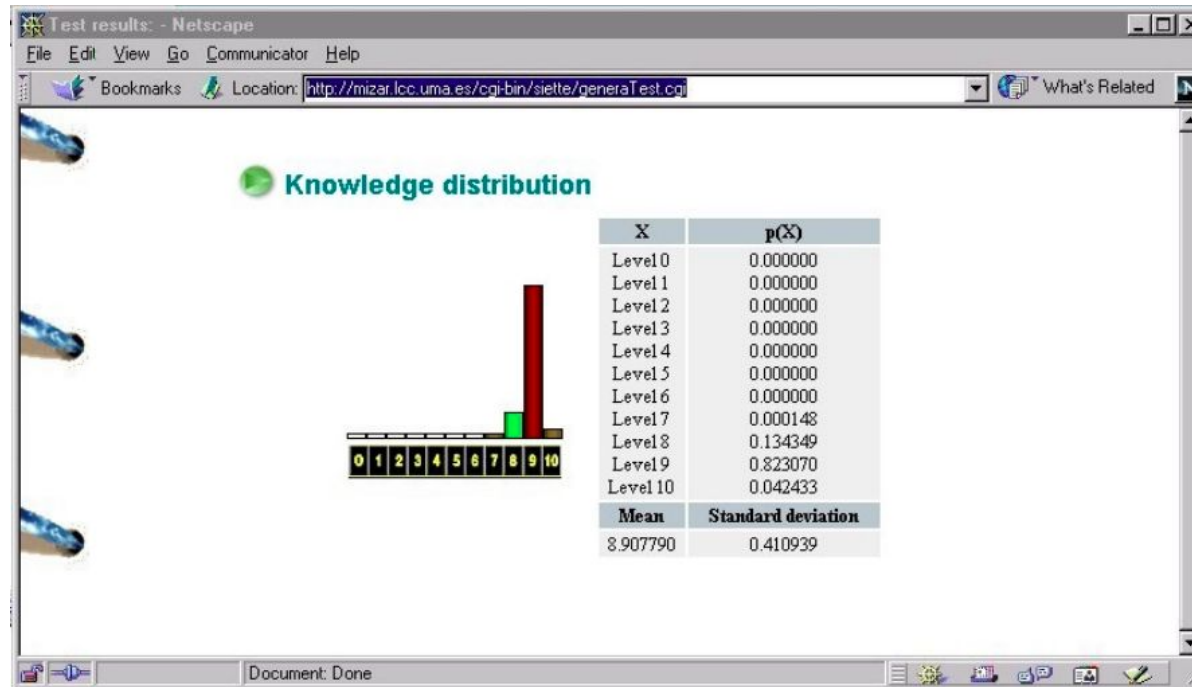


Fig. 6. Student state after finishing a test. (Conejo, et al. 2004)

- Evaluation with simulated students
  - Number of classes  $K$  and confidence factor  $\rho$
  - Discrimination  $a$  and guessing factor  $c$
  - Item selection method (Bayesian vs. difficulty-based vs. random)
  
- Evaluation with real students

- Setup:  $N$  randomly generated students for each discrete ability level in  $[0, K-1]$
- Assumption: Items are correctly calibrated
- Test Termination: Probability of student belonging to a particular ability level reaches  $\rho$

- Evaluated: Number of classes  $K$  and confidence factor  $\rho$
- Results:
  - Impossible to classify all test takers correctly
  - Percentage of correctly classified test takers dependent on  $\rho$
  - Number of questions in the test strongly correlated to  $K$

# Experiment 1



	<i>Confidence factor <math>\rho = 0.75</math></i>		<i>Confidence factor <math>\rho = 0.90</math></i>		<i>Confidence factor <math>\rho = 0.99</math></i>	
<i>Number of classes <math>K</math></i>	<i>% of correctly classified students</i>	<i>Average number of questions posed <math>T</math></i>	<i>% of correctly classified students</i>	<i>Average number of questions posed <math>T</math></i>	<i>% of correctly classified students</i>	<i>Average number of questions posed <math>T</math></i>
3	84.05	2.00	95.82	3.58	99.46	5.65
5	81.61	6.23	92.76	10.38	99.37	19.27
7	80.96	11.11	92.85	18.16	99.38	33.12
9	80.86	16.15	92.93	26.39	99.42	47.27
11	80.52	21.19	92.92	34.54	99.26	60.85

Table 1. Experiment 1 results. (Conejo, et al. 2004)

- Evaluated: Discrimination  $a$  and guessing factor  $c$
- Results:
  - Best results from items with higher  $a$  and lower  $c$
  - $a$  does not influence number of questions past a certain point
  - Extremely high number of questions needed for lower  $a$



# Experiment 2

<i>Guessing factor <math>c</math></i>	<i>% of correctly classified students</i>	<i>Average number of questions posed <math>T</math></i>
0.00	92.85	18.16
0.10	92.37	25.34
0.25	92.11	36.05
0.33	91.73	43.37
0.50	91.49	63.37

<i>Discrimination index <math>a</math></i>	<i>% of correctly classified students</i>	<i>Average number of questions posed</i>
0.20	90.4	174.9
0.50	91.5	35.2
0.70	91.9	26.3
1.20	92.8	18.1
1.70	93.8	15.3
2.20	95.4	14.8

Table 2. Experiment 2 results. (Conejo, et al. 2004)

- Evaluated: Item selection method (Bayesian vs. difficulty-based vs. random)
- Results:
  - Adaptive methods (Bayesian and difficulty-based) performed better than random selection
  - Higher value of  $K$  showed even better performance of adaptive methods
  - Ultimately difficulty-based selection was chosen because of lower computational cost

# Experiment 3

	<i>Bayesian</i>		<i>Difficulty-based</i>		<i>Random</i>	
<i>Number of classes K</i>	<i>% of correctly classified students</i>	<i>Average number of questions posed T</i>	<i>% of correctly classified students</i>	<i>Average number of questions posed T</i>	<i>% of correctly classified students</i>	<i>Average number of questions posed T</i>
3	96.06	3.58	95.62	3.58	95.82	3.58
5	93.31	6.87	94.67	7.37	92.76	10.38
7	92.75	8.70	94.43	9.03	92.85	18.16
9	92.53	9.85	94.23	10.14	92.93	26.39
11	92.10	10.71	94.14	11.02	92.92	34.54

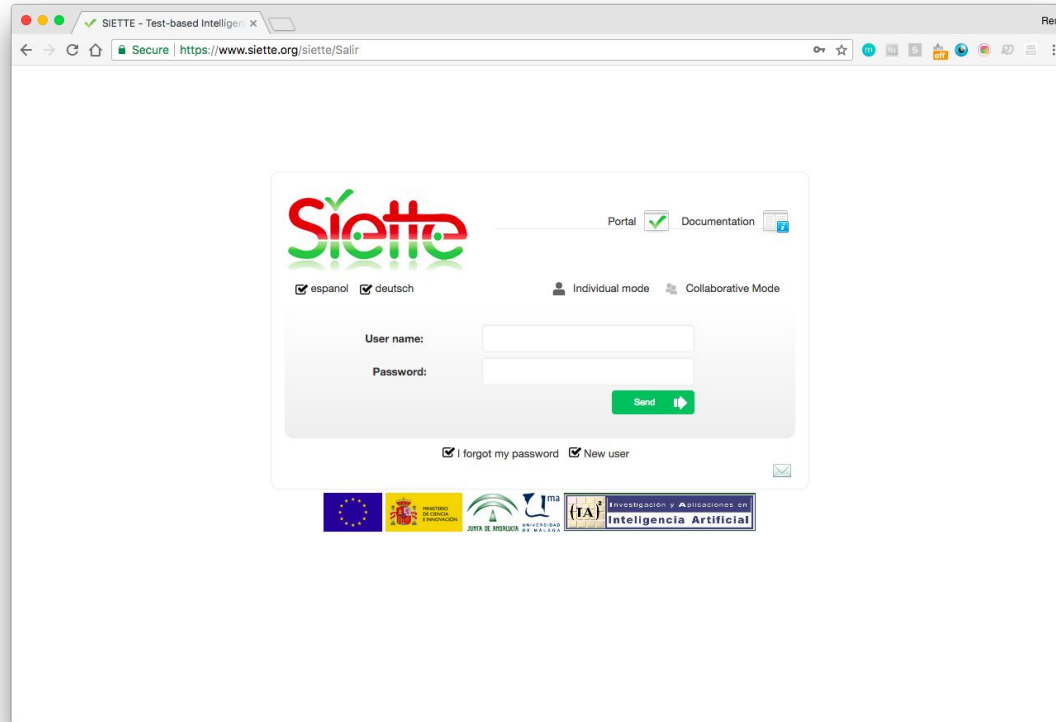
Table 3. Experiment 3 results. (Conejo, et al. 2004)

- Informal evaluation study in June 2000
- 24 Computer Science students of the subject Artificial Intelligence at the University of Malaga, Spain
- Students had an exam in about four weeks
- Students took an online test about LISP then filled out a questionnaire
- Participation was voluntary and anonymous

- Majority needed less than five minutes to use the system
- Majority considered SIETTE useful
- Half said that they would use the tool again, the rest were either unsure or said that they would not
- Majority would recommend it to other students

- Majority considered the grade they received right, fair, or totally fair
- Majority considered difficulty of the test normal
- Tie between pen-and-paper test and SIETTE test, mainly because they would rather wait for SIETTE to be fully developed and tested

# Updates on SIETTE



- SIETTE can be integrated with Moodle LMS
- No longer just a test-based system but an automatic assessment environment
- Includes more question forms such as short answer questions
- Includes other assessment methods including traditional percentage scoring
- Only 18 tests used IRT, compared to 642 which used percentage scoring and 581 which used item scoring



# Issues and Challenges

- General resistance to change
- Difficulty in understanding underlying concepts
- Item calibration

# Conclusions

- Item Response Theory and Computer Adaptive Tests are a valid alternative to Classical Test Theory and traditional pen-and-paper testing.
- The use of IRT and CAT can reduce test length and therefore test time.
- Members of every level of the education system have to be educated about IRT and CAT before its use is more widely accepted.
- More effective and efficient ways of item difficulty estimation and item calibration should be studied and implemented.

# References

Baker, Frank B. *The Basics of Item Response Theory*. 2nd ed., ERIC Clearinghouse on Assessment and Evaluation, 2001, 1-57.

Barla, Michal, et al. "On the impact of adaptive test question selection for learning efficiency." *Computers & Education* 55.2 (2010): 846-857.

Conejo, Ricardo, et al. "SIETTE: A web-based tool for adaptive testing." *International Journal of Artificial Intelligence in Education* 14.1 (2004): 29-61.

Conejo, Ricardo, Eduardo Guzmán, and Monica Trella. "The SIETTE automatic assessment environment." *International Journal of Artificial Intelligence in Education* 26.1 (2016): 270-292.

Thissen, David. "Bad questions: An essay involving item response theory." *Journal of Educational and Behavioral Statistics* 41.1 (2016): 81-89.

Wauters, Kelly, Piet Desmet, and Wim Van Den Noortgate. "Adaptive item-based learning environments based on the item response theory: possibilities and challenges." *Journal of Computer Assisted Learning* 26.6 (2010): 549-562.

Thank You!