

IMPERIAL COLLEGE LONDON

BENG INDIVIDUAL PROJECT - FINAL REPORT

**jSCAPE - Java Self-assessment
Center of Adaptive Programming
Exercises**

Author:
Alexis CHANTREAU

Supervisor:
Dr. Tristan ALLWOOD

June 6, 2014

Contents

1	Introduction	2
1.1	Motivation	2
1.2	Objectives	2
1.3	Contributions	4
1.4	Report Structure	4
2	Background	5
2.1	Computer Based Tests	5
2.2	Computerized Adaptive Testing	5
2.3	Probabilistic Test Theory	9
2.3.1	Likelihood and Maximum Likelihood Estimation	9
2.3.2	Bayes probability theory	10
2.3.3	Item Response Theory	10
2.4	Summary	11
3	Related Work	12
3.1	Environment for Learning to Program	12
3.2	CourseMarker	12
3.3	Automatic Exercise Generator with Tagged Documents based on the Intelligence of Students	13
3.4	Adaptive Self-Assessment Master	13
3.5	System of Intelligent Evaluation Using Tests for Tele-education	13
4	Evaluation	14
4.1	Qualitative	14
4.2	Quantitative	14
5	Conclusion	15
5.1	Future Work	15
A	User Manual	17

List of Figures

2.1	Flowchart of an adaptive test. Adapted from [1].	7
-----	--	---

Chapter 1

Introduction

Mention three tier architecture

1.1 Motivation

- The rise in MOOCs such as Coursera, Udacity shows that there is a real interest in learning programming.
- Programming is considered "difficult" to learn and it can only be learnt effectively through lots and lots of practice.
- This is even more apparent when students are introduced to other programming languages with different paradigms such as Haskell and Prolog, or more low level programming languages such as C.
- Hence, the need to provide a platform for students to practice their programming skills and understanding of programming concepts.
- Having a lecturer come up with all the exercises by himself can be both time consuming and ineffective: some exercises may not be challenging enough for certain top students, or on the contrary too difficult for struggling students, which can be discouraging.
- Lecturers can only rely on a few homework assignments to get an idea of how students are doing. Having a way for lecturers to gather large amounts of data about students' performances would be beneficial. Allows for supplementary material, exercises, etc...

1.2 Objectives

Having identified the problems associated with teaching and learning programming, we were lead to formulating objectives in order to make the

project successful and useful to the parties involved.

The main objective of the project was to produce a web-based teaching infrastructure to complement the introductory first year programming classes. Four key features were identified:

- **Programming questions/exercises** - The web platform should allow students to practice their programming skills and understanding of programming concepts. There should be no limit to the number of questions a student can answer, so that if a student desires more practice, then he should be able to do that. Additionally, it should be possible for a specific set of people, such as lecturers and tutors, to add questions/exercises to the system.
- **Progress tracking** - Designated people, such as lecturers and tutors, should have access to detailed statistics about the students performances. This will provide them with useful information about difficulties particular students, or the entire class, may be facing. In addition, the system should give feedback to the students, in the form of simple statistics, allowing them to identify their weak areas and thus improve on them.
- **Adaptive difficulty** - The questions or exercises presented to the students should be suited to their ability. Not only will this stimulate the learning process, but it will also give a better indication of a student's understanding of the programming concepts being tested.
- **Automated generation** - There should be a mechanism to allow for some degree of automated or semi-automated generation of exercises. This will provide a large supply of "fresh" questions, so that students don't end up answering the same questions and memorizing the answers to them.

While investigating existing solutions (Section 2.4 - Related Work) we found out that some of these features were less common than others. The availability of programming exercises and progress tracking are very essential in such systems, therefore all of the tools implemented these features. On the other hand, relatively few tools integrated some form of adapting the questions to the students' ability. Finally, almost none of the tools featured automated generation of questions, opting instead to allow exercises to be added manually to the system.

1.3 Contributions

1.4 Report Structure

Chapter 2

Background

2.1 Computer Based Tests

CBT abbreviation - offers advantages such as being able to display higher quality visuals such as pictures, videos, graphs, etc... - low paperwork, everything is stored on the computer - automatic grading, less work for teachers - generation of statistics is made easier by the fact that the data can be processed by the computer - nowadays a lot of learning is done on computer systems, and children are used to dealing with computers so assessment through this medium is advantageous.

CBTs are typically "fixed-item" tests where all the students answer the same set of questions, usually provided by the person responsible for the assessment. This isn't ideal since students can be presented with questions that are too easy or too difficult for them to answer. Consequently, the results of the test won't give a very accurate representation of a student's ability, and for this reason, these types of tests aren't extremely useful. This problem lead to research and the development of computerized adaptive testing (CAT).

2.2 Computerized Adaptive Testing

Computerized adaptive testing (CAT), also called *tailored testing*, is a form of computer-based testing which administers questions (referred to as *items* in the psychometrics literature) of the appropriate difficulty by adapting to the examinee's ability. For example, if an examinee answers an item correctly, then the next item presented will higher on the difficulty scale. On the other hand, if they answer incorrectly, they will be presented with an item lower on the difficulty scale.

From an architectural perspective, a computerized adaptive test (CAT) consists of five components [3]:

1. Calibrated item pool

An item pool is needed to store all the items available for inclusion in a test. This item pool must be calibrated with a psychometric model. During this phase, the item parameters are estimated according to the chosen model and scaled to fit with already existing items. Usually, the psychometric model employed in these systems is called Item Response Theory (IRT) (section 2.3.3). Calibration is a complex process, and to be done accurately it requires a considerable amount of data. Typically, it is performed by psychometricians, aided by expensive and sophisticated calibration software.

2. Starting point

Initially, when zero items have been administered, no information is known about the examinees and so the CAT is unable to estimate their ability. As a result, the item selection algorithm will fail to choose the next item to be administered. If there is previous information available, for example an examinee's ability estimate in a closely related subject, then this can be input into the system to form the starting point configuration. Often this data isn't available or too costly to collect, so the CAT's initial ability estimate for the examinee corresponds to the mean on the ability scale - hence the first item presented will be of average difficulty.

3. Item selection algorithm

The item selection algorithm chooses the next item to present to the examinee based on the ability estimate of the examinee up to that point. Several methods exist and largely depend on the psychometric model in use. One of the most commonly used methods is the *maximum information method* (section 2.3.3), which selects the item which maximizes the information function with respect to the estimated ability at that point.

4. Scoring algorithm

The scoring algorithm refers to the steps taken to update the examinee's ability estimate after an item has been answered. The two most commonly used methods are *maximum likelihood estimation* (section 2.3.3) and *Bayesian estimation* (section 2.3.3).

5. Termination criterion

The termination criterion specifies when the CAT should finish. For example the CAT can terminate when the change in the ability estimate after each

iteration is below a certain threshold, or when time has run out, or when N items have been administered, etc... Obviously, the CAT shouldn't be terminated too early, so as to allow enough time to estimate the examinee's ability with acceptable accuracy.

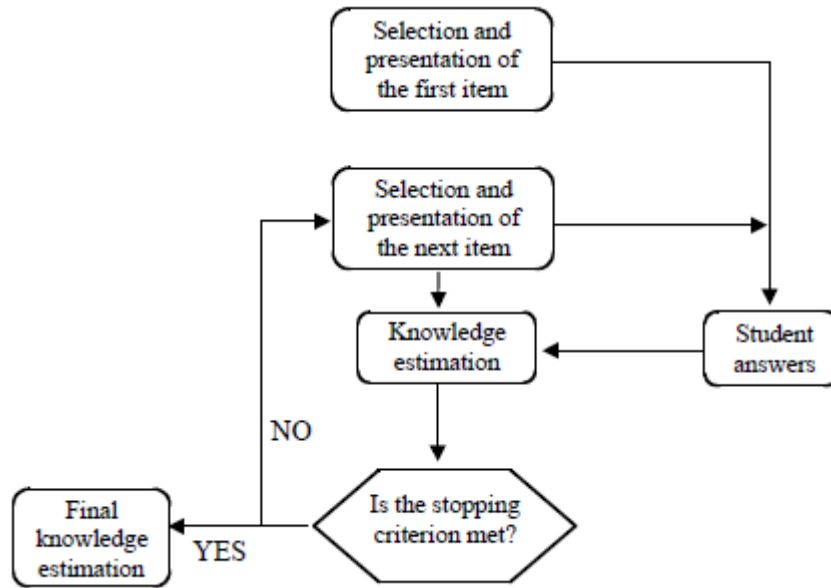


Figure 2.1: Flowchart of an adaptive test. Adapted from [1].

The flowchart in figure 2.1 corresponds to components 2-5, and illustrates the basics of the algorithm implemented in CAT. [2] gives a more detailed description of the procedure:

1. The pool of items that haven't been administered yet is searched to determine the best item to present to the examinee, according to the current estimation of his ability.
2. The chosen item is presented to the examinee, who then answers it correctly or incorrectly.
3. The ability estimate is updated, based upon this new piece of information and the previous ability estimate.
4. Steps 1–3 are repeated until a termination criterion is met.
5. The algorithm returns a final ability estimate for the examinee's performance along with a confidence level: a percentage value indicating how accurate the estimate is.

CATs offer several advantages over traditional CBTs. As a result CATs have been used in many areas[4], such as education, job hiring, counselling, clinical studies, etc... Since CATs administer items by adapting to the examinee's ability, the test-taking experience ends up being a more positive one. Indeed, examinees won't have to deal with answering items which are too difficult or too easy compared to their ability level, a problem which appears in traditional CBTs.

In addition, by administering only those items which will yield additional information, CATs end up being more accurate in estimating an examinee's ability level. This contrasts with CBTs which usually provide the best precision for examinees of medium ability, whereas extreme scores end up being less accurate.

Lastly, CATs can come up with an ability estimate much quicker and with fewer administered items when compared to traditional CBTs. Indeed, an adaptive test can typically be shortened by 50% and still maintain a higher level of precision than a fixed version.[5]

Despite the advantages mentioned above, CATs have some limitations. A frequent complaint is that an examinee isn't allowed to go back and change his answer to a past item. This limitation exists to prevent the examinee from intentionally answering items incorrectly to make subsequent items easier, and then going back and selecting the correct answers to achieve a perfect score. For similar reasons, it isn't possible to skip items, the examinee must select an answer to move on to the next item.

The second issue has to do with the items themselves. First of all, there is the need for a large bank of items to cater to all ability levels. Developing an item pool of sufficient size can be very time consuming. David J. Weiss writes in [6] that item pools with 150-200 items are to be preferred, although 100 high quality items can sometimes be enough to achieve adequate estimations of ability levels.

Secondly, for the CAT to be of good quality the item pool needs to be calibrated accurately. This requires pre-administering the items to a sizeable sample and then simultaneously estimating all the item parameters for each item. The guidelines in [7] suggest that sample sizes may be as large as 1000 examinees. This phase is costly, time consuming and often times simply unfeasible.

Lastly, item exposure is a possible security concern. Sometimes particular items may be presented too often and become overused. This may result in examinees becoming familiar with them and sharing them to other examinees of the same ability level, thus corrupting the results of the test.

This problem can be solved to some extent by modifying the item selection algorithm to include some exposure control mechanism.

A brief overview of CATs was given in this section. All of these concepts will be explored in more detail in item response theory (section 2.3.3) and in the implementation of adaptive testing in jSCAPE (chapter XX?).

2.3 Probabilistic Test Theory

In this section we go over a few topics in probability and how they can be implemented successfully in the area of computer assessments.

2.3.1 Likelihood and Maximum Likelihood Estimation

Although the terms probability and likelihood are used interchangeably in every day life, in statistics a distinction can be made.

For any stochastic process, let us denote the observed outcomes as x and the set of parameters as θ . When we say probability, we want to calculate $P(x|\theta)$, i.e. the probability of observing the outcomes x given specific values for the set of parameters θ .

In more mathematical terms, we have:

$$L(\theta|x) = P(x|\theta)$$

To highlight the distinction we illustrate with an example of how the two terms are used. If we consider a dice, a possible parameter is the fairness of the dice, while possible outcomes are which values are displayed after a roll. For instance, if a fair dice is rolled 5 times, what is the *probability* that a 6 will show up on every roll? If a dice is rolled 5 times and lands on 6 every roll, what is the *likelihood* that the dice is fair?

Maximum likelihood estimation refers to a method of statistical inference where one can find the set of parameter values (θ) which are most *likely* given the observed outcomes (x). As mentioned in the name of the method, this is done by finding the parameter values which maximize the likelihood function $L(\theta|x)$.

2.3.2 Bayes probability theory

2.3.3 Item Response Theory

We mentioned in section 2.2 that Item Response Theory (IRT) is usually the psychometric model of choice when developing a CAT, e.g. the Graduate Record Examination (GRE) and Graduate Management Admission

Test (GMAT). CATs can still be implemented with Classical Test Theory but they offer less sophistication and less information to evaluate/improve the reliability of the test, making IRT the superior choice.

IRT attempts to model the answer to an item as a mathematical model.

Several IRT models have been developed over the years to address the different types of tests that exist, e.g. multiple choice exams, agreement questionnaires (Likert scale), etc... These models all seek to achieve the same goal, modelling the examinee's ability on some ability scale, but they differ in the number of parameters associated with each item. We now take a greater look at these different models:

The one-parameter logistic model

$$P(\theta) = \frac{1}{1 + e^{-(\theta-b)}} \quad (2.1)$$

Equation (2.1) shows the item response function for the one-parameter logistic model (1PL).

The two-parameter logistic model

$$P(\theta) = \frac{1}{1 + e^{-1.7a(\theta-b)}} \quad (2.2)$$

Equation (2.2) shows the item response function for the two-parameter logistic model (2PL).

The three-parameter logistic model

$$P(\theta) = c + \frac{1 - c}{1 + e^{-1.7a(\theta-b)}} \quad (2.3)$$

Equation (2.3) shows the item response function for the three-parameter logistic model (3PL).

Estimating the ability

In section 2.2, when discussing scoring algorithms, we saw that the two most common methods for ability estimation were *maximum likelihood estimation* and *Bayesian estimation*.

In the *maximum likelihood estimation* method, we need to define a likelihood function in terms of the ability level (θ) we are trying to estimate:

$$L(\theta|\mathbf{u}) = L(\theta|u_1, \dots, u_n) = P(u_1, \dots, u_n|\theta) = \prod_{i=1}^n P_i(\theta)^{u_i} (1 - P_i(\theta))^{(1-u_i)}, \quad (2.4)$$

where $\mathbf{u} = (u_1, \dots, u_n)$ is called the response vector for an examinee, that is $u_i = 1$ if the examinee answers the i^{th} item correctly, and $u_i = 0$ if the examinee answers the i^{th} item incorrectly. $P_i(\theta)$ corresponds to the probability of answering the i^{th} item correctly when the ability level is θ , and thus $1 - P_i(\theta)$ gives the probability of answering the i^{th} item incorrectly when the ability level is θ .

Now that the likelihood function is defined, we can apply maximum likelihood estimation to find the ability level which is most likely given the examinee's response vector. Let us illustrate with a very simplistic example where ability levels are discrete values between -1 and $+1$. We would iterate through these ability levels and compute the following values, for example:

$$\begin{aligned} L(\theta = -1|\mathbf{u}) &= 0.001 \\ L(\theta = +0|\mathbf{u}) &= 0.012 \\ L(\theta = +1|\mathbf{u}) &= 0.054 \end{aligned}$$

The likelihood function is maximized when $\theta = +1$, and so this is the maximum likelihood estimate for θ . In a CAT, and based on the examinee's response vector \mathbf{u} , the system would assign $+1$ as the ability level for that examinee.

The *Bayesian estimation* method is also somewhat based on the likelihood function and therefore very similar to *maximum likelihood estimation*.....

Item information

2.4 Summary

We have gone over types of computer based assessment and how they differ from regular assessment techniques. We have also provided an overview of the theoretical components and probability concepts which will be implemented in jSCAPE (chapter xx?).

Chapter 3

Related Work

Web-based/Computer based education and adaptive web-based assessment systems are a “hot” research area, and as a result, numerous tools, environments and infrastructures have emerged. There are common features to all, however some distinguish themselves by having not so common features. Automated exercise generation in these tools is usually non-existent or very limited. Moreover, the tools are more focused on assessing students rather than self-assessment, i.e. students take tests which count towards their final grade on these systems.

In this section we look at related software and evaluate them with respect to the objectives listed at the beginning of the development of jSCAPE.

3.1 Environment for Learning to Program

Environment for Learning to Program (ELP) is an interactive web based environment for teaching programming to first year Information Technology students at Queensland University of Technology (QUT).

3.2 CourseMarker

CourseMarker is a re-design of Ceilidh, a computer based assessment system used at the University of Nottingham for 13 years. Ceilidh was quite a complete system, providing coursework, the management of modules and the presentation of module content.

3.3 Automatic Exercise Generator with Tagged Documents based on the Intelligence of Students

The Automatic Exercise Generator with Tagged Documents based on the Intelligence of Students (AEGIS)

3.4 Adaptive Self-Assessment Master

Adaptive Self-Assessment Master (ASAM) is an extension to CourseMarker, which improves upon it by administering questions which are suited to the student's ability.

3.5 System of Intelligent Evaluation Using Tests for Tele-education

The System of Intelligent Evaluation Using Tests for Tele-education (SI-ETTE) is a web based environment for generating and constructing adaptive tests.

Chapter 4

Evaluation

Mention how our developed system performs against the advantages and disadvantages of CAT.

The evaluation stage will address mostly these two aspects:

- The system has correctly modelled the ability of students
- The system is useful in helping students to learn programming and helping lecturers with getting feedback on their teaching, in the form of statistics.

4.1 Qualitative

- Surveys to get feedback from students on interface, usability, etc...
-

4.2 Quantitative

- Statistical analysis to evaluate item calibration and modelling of student's abilities
-

Chapter 5

Conclusion

5.1 Future Work

- Very flexible system so other programming languages could be offered, i.e. cSCAPE, for C and hSCAPE for Haskell.
- Working more extensively on the adaptive component of the system, i.e. improving the algorithm which selects questions for students based on their estimated ability.
- Extend system to allow admins to provide their own question templates, maybe come up with a template grammar which then allows questions to be automatically generated. Or at least allow pluggable function references which will be called to generate the exercise component.
- Add support for more question types. JavaFX is very good in that sense since it can play audio clips, video clips, show animations, the webview component has endless possibilities thanks to the inclusion of Javascript.
- Future Work as a research project vs future work as a commercial product
- JavaFX applications are based on the model-view-controller pattern, so a nice split can be done in the code, however I only learnt about this two weeks into the project, so all the of the GUI components are created in the code as opposed to in the FXML file.

Bibliography

- [1] Conejo, R., Guzmán, E., Millán, E., Trella, M., Pérez-de-la-Cruz, J. L., & Rios, A. (2004). SIETTE: A Web-Based Tool for Adaptive Testing. *International Journal of Artificial Intelligence in Education*, 14, 29-61.
- [2] Computerized adaptive testing. http://en.wikipedia.org/wiki/Computerized_adaptive_testing. Accessed: June 6, 2014.
- [3] Thompson, Nathan A., & Weiss, David A. (2011). A Framework for the Development of Computerized Adaptive Tests. *Practical Assessment, Research & Evaluation*, 16(1).
- [4] IRT-Based CAT. <http://www.iacat.org/content/irt-based-cat>. Accessed: June 6, 2014.
- [5] Weiss, D. J., & Kingsbury, G. G. (1984). Application of computerized adaptive testing to educational problems. *Journal of Educational Measurement*, 21, 361-375
- [6] Weiss, D.J. (1985). Adaptive Testing by Computer, *Journal of Consulting and Clinical Psychology*. 1985, 53, 6, pp. 774-789.
- [7] Wainer, H., & Mislevy, R.J. (2000). Item response theory, calibration, and estimation. In Wainer, H. (Ed.) *Computerized Adaptive Testing: A Primer*. Mahwah, NJ: Lawrence Erlbaum Associates.

Appendix A

User Manual