

Educational Assessment Using Intelligent Systems

*Valerie J. Shute
Diego Zapata-Rivera*

December 2008

ETS RR-08-68



Educational Assessment Using Intelligent Systems

Valerie J. Shute¹

Florida State University, Tallahassee

Diego Zapata-Rivera

ETS, Princeton, NJ

December 2008

As part of its nonprofit mission, ETS conducts and disseminates the results of research to advance quality and equity in education and assessment for the benefit of ETS's constituents and the field.

ETS Research Reports provide preliminary and limited dissemination of ETS research prior to publication. To obtain a PDF or a print copy of a report, please visit:

<http://www.ets.org/research/contact.html>

Copyright © 2008 by Educational Testing Service. All rights reserved.

ETS, the ETS logo, and LISTENING. LEARNING.
LEADING. are registered trademarks of Educational Testing
Service (ETS).



Abstract

Recent advances in educational assessment, cognitive science, and artificial intelligence have made it possible to integrate valid assessment and instruction in the form of modern computer-based intelligent systems. These intelligent systems leverage assessment information that is gathered from various sources (e.g., summative and formative). This paper analyzes the role of educational assessment in intelligent systems, summarizes the characteristics of successfully deployed intelligent systems, and describes an evidence-based approach to incorporating valid and reliable assessments into enhanced intelligent systems.

Key words: Evidence-centered design, formative assessment, summative assessment, intelligent tutoring system, reliability, student model, validity

Acknowledgments

We'd like to acknowledge Eric Hansen, Irvin Katz, Don Powers, and Bob Linn for their sage comments on an earlier draft of this paper. Finally, a version of this paper appears as an entry in the *Third Edition of the International Encyclopedia of Education*, edited by E. Baker, B. McGaw, & P. Peterson.

Introduction

Models of educational measurement influence instructional practices in the classroom and thus have different effects on student learning. The primary goal of measuring students' educational progress historically has been to identify differences in achievement among students to rank order them. The types of measurement models used to achieve this goal rely on summative assessment, which is useful for accountability purposes but not as helpful for guiding day-to-day instruction. Student-centered measurement models, in contrast, rely on formative assessment, which are useful in guiding instruction and supporting student learning but does not provide enough of a basis for accountability purposes. Now, however, it may be possible to combine these models within new, enhanced intelligent systems. These systems can use both kinds of assessments—summative and formative—and harness computer technology, educational measurement, and cognitive science to address problems in education.

This unified approach to educational measurement rests on the following assumptions:

1. Individual differences among students affects learning.
2. Such effects can be quantified and predicted.
3. Technology can capitalize on these effects to benefit primarily teachers and students, but also others involved in the educational process, such as administrators and parents.

The goal is to figure out how to integrate assessment and instruction to improve student learning and education.

This paper defines educational measurement, in terms of the role assessment plays in education, and intelligent systems focusing on computer usage of assessment data to make inferences about students' cognitive and other attributes. The paper also examines the role that assessment plays in both traditional and enhanced intelligent systems and provides an outline of an approach to incorporating evidence-based assessments into intelligent systems to improve learning.

Definitions

Educational Measurement

Educational measurement may be defined broadly as the application of a standard scale or measuring tool to determine the degree to which educationally valuable knowledge, skills, and abilities have been acquired. According to the National Council on Measurement in Education

(2008), this includes theory, techniques, and instrumentation available for measurement of educationally relevant human, institutional, and social characteristics. We measure to obtain information, and such information may or may not be useful, depending on the accuracy of the instruments and on the skillful manner with which they are used.

Assessment is a general term that includes testing. Progress toward educational goals is typically assessed through testing. Assessment is both an instrument and a process by which information is obtained relative to a known objective or goal. Since inferences are made about what a person knows on the basis of or her responses to a limited number of assessment tasks or items, there is always some uncertainty in inferences made on the basis of assessments. The goal in educational measurement is to minimize uncertainty or error; thus, key aspects of assessment quality are validity and reliability. Reliability refers to the consistency of assessment results—the degree to which they rank order students in the same way. Validity refers to the extent to which the assessment accurately measures what it is supposed to measure and the accuracy of the inferences made from task or test results.

Types of Assessment

Here are considered two main types of assessment: summative and formative. Summative assessment reflects the traditional approach used to assess educational outcomes. This involves using assessment information for high-stakes, cumulative purposes, such as promotion, certification, and so on. Summative assessment is usually administered after some major event, like the end of the school year or marking period, or before a big event, like college. Three benefits of this approach are the following: (a) It allows for comparing student performances across diverse populations on clearly defined educational objectives and standards, (b) it provides reliable data (e.g., scores) that can be used for accountability purposes at various levels (e.g., classroom, school, district, state, and national) and for various stakeholders (e.g., students, teachers, and administrators), and (c) it can inform educational policy (e.g., curriculum or funding decisions).

Formative assessment reflects a more progressive approach in education. This involves using assessments to support teaching and learning. Formative assessment is tied directly into the fabric of the classroom and uses results from students' activities as the basis on which to adjust instruction to promote learning in a timely manner. This type of assessment is administered much more frequently than summative assessment and has shown great potential for harnessing

the power of assessments to support learning in different content areas and for diverse audiences. In addition to providing teachers with evidence about how their students are learning so that they can revise instruction appropriately, formative assessment may directly involve students in the learning process, such as by providing feedback that will help students gain insight about how to improve and by suggesting (or implementing) instructional adjustments based on assessment results (Black & Wiliam, 1998a, 1998b; Shute, 2007).

We now turn our attention to intelligent, computer-based systems that have been around for several decades but have yet to be fully embraced by education. Their primary goal is to enhance student learning, so assessment, in theory, should play a key role in these systems.

Intelligent Systems

Intelligent systems (also known as intelligent tutoring systems) refer to educational software containing an artificial intelligence component. The software tracks students' work, adjusting feedback and providing hints along the way. By collecting information on a particular student's performance as well as other cognitive and noncognitive variables, the software can make inferences about strengths and weaknesses and can suggest additional work.

A summary of requirements for intelligent systems was presented by Hartley and Sleeman (1973). They argued that these systems must possess (a) knowledge of the learner (student model), (b) knowledge of the domain (expert model), and (c) knowledge of teaching strategies (pedagogical model). It is interesting to note that this simple list has not changed in more than three decades; however, advances have been made in each of the three areas. All of the computer-resident knowledge marks a radical shift from earlier knowledge-free, computer-assisted instructional programs. Furthermore, the ability to diagnose students' errors and adapt instruction based on the diagnosis represents a key difference between intelligent and other computer-based systems, such as simulations. Intelligent systems are also aligned with the features and goals of formative assessment. The three main components of intelligent systems—student, expert, and pedagogical models—are now briefly described.

A student learns from an intelligent system primarily by solving problems—ones that are appropriately selected or tailor made and that serve as learning experiences for that student. The system may start by assessing what the student already knows. Information about the student is maintained within what is called the student model, which is updated during the course of learning. The system then must consider what the student needs to know. This information is

embodied in the domain expert model. Finally, the system must decide what unit of content (e.g., assessment task or instructional element) ought to be presented next, and how it should be presented. This is achieved by the pedagogical model (or tutor). From all of these considerations, the system selects or generates a problem and then either works out a solution to the problem (via the domain expert model) or retrieves a prepared solution. The intelligent system compares its solution to the one the student has prepared and performs a diagnosis based on differences between the two as well as on other information available in the student model. Feedback is offered by the system based on considerations, such as how long it has been since feedback was last provided, whether the student already received some particular advice, and so on. After this, the program updates the student model, and the entire cycle is repeated, starting with selecting or generating a new problem.

Despite the great promises of intelligent systems, they are currently not widely used in classrooms, partly because of their cost and also because of measurement limitations. We now focus on the latter in more detail, describing how assessments differ between traditional intelligent systems and newer, enhanced intelligent systems. This is intended to provide the foundation on which to consider a new view of educational measurement within intelligent systems.

Assessments' Role in Intelligent Systems

For the most part, traditional intelligent systems use a formative assessment model, where different student actions invoke different instructional decisions or paths. This comprises the basis for adaptive instruction. New, enhanced intelligent systems extend the assessment capabilities of traditional systems. Some of these enhancements include the use of evidence-based assessment data, explicit links to state curriculum standards, formative and summative sources of assessment information, new measurement techniques from educational psychology and cognitive science, and an explicit and strong role for teachers. Both types of intelligent systems are now discussed in turn.

Traditional Intelligent Systems

As noted earlier, formative assessment is explicitly intended to support student learning, defining the role of the student as an active, creative, and reflective participant in the learning process. Learning environments that make use of formative assessment typically include

individualized instruction, along with hands-on, authentic learning activities. Assessments are used primarily to inform teaching and improve student learning.

One major downside of this model is that formative assessment is often implemented in a nonstandardized and hence less rigorous manner than summative assessment. This can hamper the validity and reliability of the assessment tools and data. The validity and reliability of the assessment data affect the accuracy of the student diagnosis, and the diagnosis informs instructional support. Therefore, if the first part of the chain is weak, the rest (i.e., diagnostic accuracy and effective instructional support) consequently would be compromised. In other words, the effectiveness of an intelligent system in achieving its goal hinges on the quality of information in the student model (i.e., the inferences about what the student knows and can do).

Traditional intelligent systems that employ formative assessment utilize a rich source of student data from which to draw inferences. For example, evidence is captured from all past and current student-system interactions and may differ in type and grain size. Thus, in addition to the nonstandardization of methods for implementing formative assessment in traditional intelligent systems, there are problems with accurately modeling student knowledge within such multifaceted environments. This poses a number of psychometric challenges (e.g., modeling of multiple abilities, capabilities, and other learner characteristics), regardless of the measurement model employed.

We now take a closer look at new intelligent systems that are starting to integrate formative and summative sources of assessment information. These new systems are employed within real classroom settings.

Enhanced Intelligent Systems

Most current intelligent systems reside primarily in the laboratory. This isolation from real classrooms explains why their designs have not been overly concerned with summative types of assessment and also explains to some extent why they have not been widely adopted. That is, learning systems deployed within laboratory-based environments do not have to comply with the same high standards (e.g., accountability requirements), as those in real classroom environments. However, as these systems move out of the laboratory and into the classroom, the need for accountability (e.g., standards and norm-referenced assessments) increases.

Summative assessments are explicitly designed for accountability purposes. They represent a source of valid and reliable evidence of student knowledge. Because of national and

international accountability requirements and interests, summative assessments are widely used in schools. For example, in the United States, summative assessments have received increased attention after the U.S. Congress passed the No Child Left Behind Act of 2001 (NCLB, 2002). And on the international front, the OECD Programme for International Student Assessment (PISA) (2008) is being used to compare student achievement in countries all over the world. The measurement community has made important advances in the development of psychometric models (e.g., Rasch, item response theory) that provide reliable and valid assessment information, typically presented as a single measure of ability at a particular point in time for any given student. These data, however, have limited use for formative purposes. One often cited downside of this emphasis on accountability is that teachers tend to view testing as time taken away from valuable instruction and learning (e.g., Kahl, 2003).

Over a decade ago, Snow and Mandinach (1991) called for the development of principles for creating valid and useful instructional-assessment systems. Only now are intelligent systems beginning to enter classrooms that integrate sound assessment and instruction. These systems are characterized by three elements: (a) a strong presence of teachers in all phases of the project, (b) a cognitive model that is used to drive instructional and assessment interactions, and (c) explicit connections to state standards and standardized state tests.

An example of a successfully deployed intelligent system can be seen in the Web-based Cognitive Tutors (Anderson, Corbett, Koedinger & Pelletier, 1995). A derivation of their cognitive tutor approach is called Assistments (Razzaq et al., 2007) —the merging of robust assessment with instructional assistance into one system. Assistments use real (i.e., released) items from the Massachusetts Comprehensive Assessment System (MCAS) state exams within the system for both assessment and instructional purposes. Table 1 summarizes the main features that separate traditional from enhanced intelligent systems with regard to the role assessments play.

While Assistments provide a good example of joining formative and summative models within an intelligent system, this important blending is still uncommon, despite calls for their union. A few other systems demonstrating similar capabilities are SIETTE (Conejo, Guzman, Millán, Trella, Perez-DeLa Cruz, et al., 2004), ACED (Shute, Hansen, & Almond, 2007), and English ABLE (Zapata-Rivera, VanWinkle, Shute, Underwood, & Bauer, 2007). The next section presents an evidence-based approach designed to create valid assessments for summative or formative purposes, which may be implemented as part of an intelligent system.

Table 1

Assessments' Role in Traditional Versus Enhanced Intelligent Systems

Issue	Traditional systems	Enhanced systems
Design methods based on evidentiary argument	Mostly absent	Present
Assessment focus	Mostly formative assessment	Formative and summative assessment
Links to standards	Mostly absent	Present
Measurement models	Largely ad hoc	Variegated, informed by advances in educational measurement and cognitive science
Evaluations	Mostly laboratory-based	Classroom-based
Role of teacher	Very limited or absent	Strong

Evidence-Centered Design and Intelligent Systems

Evidence-Centered Approach to Assessment Design

An intelligent system that includes valid assessments, for formative or summative purposes, must elicit behavior from the student that bears evidence about key skills and knowledge. In addition, the system must provide principled interpretations of that evidence in terms that suit the purpose of the assessment. Figure 1 sketches the basic structures of an evidence-centered approach to assessment design, ECD (Mislevy, Steinberg, & Almond, 2003).

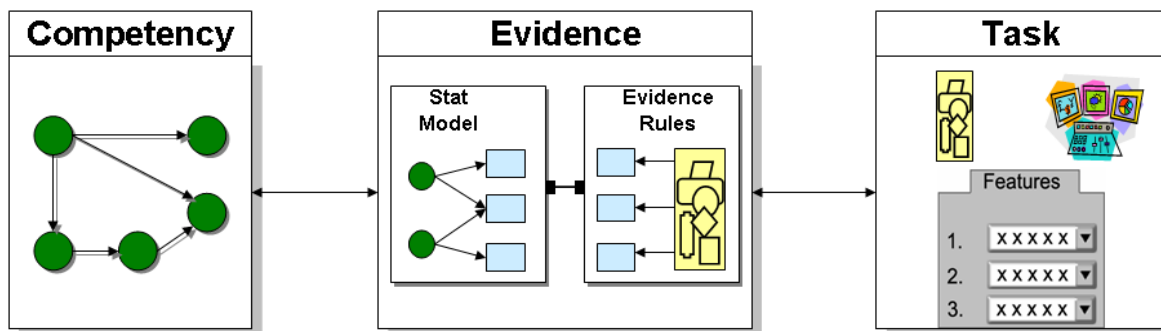


Figure 1. The three central models of an evidence-centered assessment design.

Working out these variables and models and their interrelationships is a way to answer a series of three questions posed by Messick (1992) that get at the very heart of assessment design:

1. What complex of knowledge, skills, or other attributes should be assessed? A given assessment—formative or summative—is meant to support inferences for some purpose, such as a licensing decision, provision of diagnostic feedback, guidance for further instruction, or some combination. Variables in the competency model describe the knowledge, skills, and abilities on which the inferences are to be based. The term *student model* is often used to denote a student-instantiated version of the competency model. That is, values in the student model express the assessor’s current belief about a student’s level on variables within the competency model.
2. What behaviors or performances should reveal those constructs? An evidence model expresses how the student’s interactions with and responses to a given problem constitute evidence about student-model variables. Observable variables summarize aspects of specific task performances and may come from either formative or summative sources. Then, depending on the type and origin of the sources of evidence, different parameters are used to update the student model.
3. What tasks or situations should elicit those behaviors? Task-model variables describe features of tasks or situations that will be used to elicit performance. A task model provides a framework for characterizing and constructing situations with which a student will interact to provide evidence about targeted aspects of knowledge. The task models will vary in line with the purpose of the assessment and its administration.

Within intelligent systems employing such evidence-based assessment, the student model would accumulate and represent belief about the targeted aspects of skill. These beliefs are often expressed as probability distributions for student-model variables. The way that this works in practice is that the evidence model extracts observables (e.g., scores) from student work (i.e., what the student says or does) and provides a way to aggregate scores that are then used to update the student model. In other words, the evidence model describes how evidence about a set of skills is connected to the competency-model variables using a psychometric model. Task models express situations that can evoke required evidence.

Based on the information in the student model, a viable approach to select and deliver content to the learner is needed—one that fits his or her needs at the time. This would provide context and coherence for delivering adaptive instruction, one of the main goals of intelligent systems. Following is an example of a model to support the select-and-deliver goal. It has been extended from the simpler two-process model that resides at the core of intelligent systems—diagnosis and prescription—and from a process model to support assessment (Mislevy et al., 2003).

Four-Process Adaptive Cycle

The success of any intelligent system to promote learning requires accurate diagnosis of student characteristics (e.g., algebra knowledge, troubleshooting skill, engagement). The collection of student information then can be used formatively for the prescription of optimal content, such as hints, explanations, hypertext links, practice problems, encouragement, metacognitive support, and so forth. Student information can also be used in a summative manner, such as providing reports on student achievement. The framework, described in this section, involves a four-process cycle (Shute & Zapata-Rivera, 2007) connecting the student to appropriate educational materials and other resources (e.g., learning objects, peers, applications, and pedagogical agents) through the use of a student model, shown as the small human icon at the top of Figure 2.

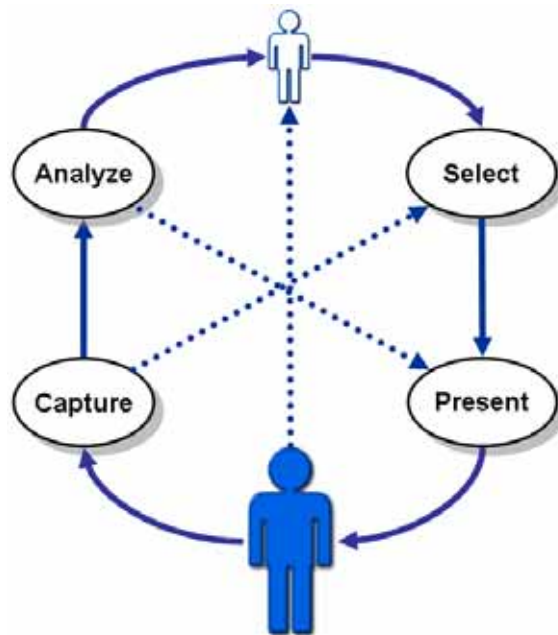


Figure 2. Four-process adaptive cycle.

The main components of this four-process cycle are (a) capture, (b) analyze, (c) select, and (d) present. The solid arrows in Figure 2 show a normal, complete loop used in many intelligent systems, whereas the dashed arrows show variations of the cycle that have been used in other kinds of systems. For example, the dashed line that goes upward from the student (represented by the large human icon in at the bottom of Figure 2) to the student model depicts an intelligent system where the student is allowed to interact directly with the student model. The nature of this interaction and the effects on the student model can vary, such as negotiating the value of a particular variable with the system or the teacher. The four processes are now briefly defined.

Analyze. The analyze process requires the creation and maintenance of a student model by properly integrating evidence sources from student performance in the environment. This usually involves representing information in the student model through inference mechanisms in relation to students' proficiency states based on specific performance data.

Select. Information (i.e., content in the broadest sense) is selected according to the model of the student maintained by the system and the goals of the system (e.g., next learning object or test item). This process is often required to determine how and when to intervene.

Present. Based on results from the select process, specific content is presented to the learner. This entails appropriate use of different media, devices, and technologies effectively and efficiently to convey information to the learner.

Discussion

Educational measurement involves making inferences about students' knowledge and skills based on limited data. How assessment data are gathered, analyzed, and used influences the student model in the following ways: (a) The granularity of assessment data affects the types of claims that can be made, such as claims about the student's general ability versus claims about the student's skills at the component level; (b) the kinds of evidence available from assessment tasks determines how reliable are the claims that can be made about the student, such as assessment claims supported using data from highly reliable test items versus claims that rely on results from homework assignments; (c) the type of evidence that can be drawn from assessment tasks helps establish the validity of the claims that can be made about the student, such as whether the student has a particular skill proficiency versus whether more evidence has to be gathered to support the claim; and (d) the sophistication of the assessment model (how assessment information is interpreted) can range from very simple to quite elaborate, such as the

use of probability-based models. It is recommended that an evidence-based assessment framework (e.g., evidence-centered design) be used to handle diverse types of assessment information—whether from summative sources or formative ones.

Intelligent systems that successfully merge instruction with information from valid formative and summative assessments could potentially improve student learning and educational outcomes. A more complete profile of the learner is needed, though, which requires new methods and new tools. These tools should (a) help create student models that provide information on aspects other than cognitive skills (e.g., noncognitive attributes such as conceptual understanding, social aspects of learning, and emotional states), (b) implement a valid assessment framework that can be used to properly analyze student work and update student models (e.g., ECD, along with the four-process adaptive cycle); (c) advise the teacher about what to do next with the class or student and how to understand the data provided by the system; and (d) encourage students to be more active and accountable for their own learning.

There is the opportunity integrate assessment and instruction into powerful new intelligent systems through recent and ongoing advances in intelligent systems, cognitive science, and educational measurement. These new systems have the potential to improve student learning and also standardized test scores. In addition, by reusing rigorous assessment tasks and items from summative tests directly in the service of learning, the chances of successful adoption of intelligent systems into mainstream education may increase.

References

- Anderson, J. R, Corbett, A. T., Koedinger, K. R., & Pelletier, R. (1995). Cognitive tutors: Lessons learned. *Journal of the Learning Sciences*, 4, 167-207.
- Black, P., & Wiliam, D. (1998a). Assessment and classroom learning. *Assessment in Education: Principles, Policy, and Practice*, 5(1), 7–74.
- Black, P., & Wiliam, D. (1998b). *Inside the black box: Raising standards through classroom assessment*. London: King’s College, School of Education.
- Conejo, R., Guzman, E., Millán, E., Trella, M., Perez-De-La-Cruz, J. L., & Rios, A. (2004). SIETTE: A Web-based tool for adaptive testing. *International Journal of Artificial Intelligence in Education* 14(1), 29-61.
- Hartley, J., & Sleeman, D. (1973). Towards more intelligent teaching systems. *International Journal of Man-Machine Studies*, 2, 215-236.
- Kahl, S. (2003, July). *Implementing NCLB assessments and accountability requirements into an imperfect world*. Paper presented at the tenth annual education law conference, Portland, ME.
- Messick, S. (1992). *The interplay of evidence and consequences in the validation of performance assessments* (ETS Research Rep. RR-92-39). Princeton, NJ: ETS.
- Mislevy, R. J., Steinberg, L. S., & Almond, R. G. (2003). On the structure of educational assessment. *Measurement: Interdisciplinary Research and Perspective*, 1(1) 3–62.
- National Council on Measurement in Education. (2008). *NCME mission*. Retrieved November 7, 2008, from <http://www.ncme.org/about/mission.cfm>
- No Child Left Behind Act of 2001, 20 U.S.C. § 6301 *et seq.* (2002).
- OECD Programme for International Student Assessment (PISA)*.. (2008) Retrieved November 7, 2008, from the Organisation for Economic Co-operation and Development: <http://www.pisa.oecd.org>
- Razzaq, L., Feng, M., Heffernan, N. T., Koedinger, K. R., Junker, B., Nuzzo-Jones, G., et al. (2007). Blending assessment and instructional assistance. In N. Nedjah, L. deMacedo Mourelle, M. Neto Borges, & N. Nunesde Almeida (Eds.), *Intelligent educational machines within the intelligent systems engineering book series* (pp. 23-49). Berlin: Springer.
- Shute, V. J. (2007). *Focus on formative feedback* (ETS Research Rep. No. RR-07-11). Princeton, NJ: ETS.

- Shute, V. J., Hansen, E. G., & Almond, R. G. (2007). *An assessment for learning system called ACED: Designing for learning effectiveness and accessibility* (ETS Research Rep. No. RR-07-26). Princeton, NJ: ETS.
- Shute, V. J., & Zapata-Rivera, D. (2007). Adaptive technologies. In J. M. Spector, D. Merrill, J. van Merriënboer, & M. Driscoll (Eds.), *Handbook of research on educational communications and technology* (3rd ed.). Mahwah, NJ: Erlbaum Associates.
- Snow, R. E., & Mandinach, E. B. (1991). *Integrating assessment and instruction: A research and development agenda*. Princeton, NJ: ETS.
- Zapata-Rivera, D., VanWinkle, W., Shute, V., Underwood, J., & Bauer, M (2007). English ABLE. In R. Luckin, K. Koedinger, & J. Greer, (Eds.), *Artificial intelligence in education: Vol. 158. Building technology rich learning contexts that work* (pp. 323–330). Fairfax, VA: IOS Press.

Further Readings

- Linn, R. L., Baker, E. L., & Dunbar, S. B. (1991). Complex, performance-based assessment: Expectations and validation criteria. *Educational Researcher*, 20(8), 15–21.
- Martin, J. D., & VanLehn, K. (1995). A Bayesian approach to cognitive assessment. In P. Nichols, S. Chipman, & R. Brennan (Eds.), *Cognitively diagnostic assessment* (pp. 141–165). Hillsdale, NJ: Erlbaum.
- Mislevy, R. J. (2006). Cognitive psychology and educational assessment. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 257–305). Westport, CT: Greenwood.
- Mislevy, R. J., & Gitomer, D. H. (1996). The role of probability-based inference in an intelligent tutoring system. *User Modeling and User-Adapted Interaction*, 5, 253–282.
- Pearl, J. (1988). *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. San Mateo, CA: Kaufmann.
- Pellegrino, J., Glaser, R., & Chudowsky, N. (Eds.). (2001). *Knowing what students know: The science and design of educational assessment*. Washington, DC: National Academy Press.
- Shute, V. J. (1995). SMART: Student modeling approach for responsive tutoring. *User Modeling and User-Adaptive Interaction*, 5(1), 1–44.
- Shute, V. J. (2007). Tensions, trends, tools, and technologies: Time for an educational sea change. In C. A. Dwyer (Ed.), *The future of assessment: Shaping teaching and learning* (pp. 139-190). New York: Lawrence Erlbaum Associates, Taylor & Francis Group.
- Shute, V. J., & Psotka, J. (1996). Intelligent tutoring systems: Past, present, and future. In D. Jonassen (Ed.), *Handbook of research for educational communications and technology* (pp. 570-600). New York: Macmillan.
- VanLehn, K. (2006). The behavior of tutoring systems. *International Journal of Artificial Intelligence in Education*, 16(3), 227-265.

Notes

¹ This paper was written while Valerie Shute was a staff member at ETS.