

Using Jitter and Sampling Techniques to Improve the Comprehensibility of Scatter Plots: A Practical Example

Michael D. Kickmeier-Rust

University of Teacher Education, St. Gallen, Switzerland

michael.kickmeier@phsg.ch

ABSTRACT: Displaying complex data including the interrelationships of several variables is one of the key challenges for information visualization. This is particularly true when the target audience has little data literacy, which is oftentimes the case in the context of learning analytics (where stakeholders are students, parents, teachers, or administrators). In this paper, I introduce a practical scenario in the context of the Swiss educational system and present an innovative solution to display complex learning data with scatter plots. By techniques such as jitter and data sampling, the scatter plot can be advanced and presented in a more comprehensible way, even when large data sets are displayed.

Keywords: Learning Performance, Performance Comparison, Visualization, Scatter Plot, Jitter

1 INTRODUCTION: SWISS TEST AND TRAINING PLATFORMS

The new Swiss national curriculum *Lehrplan 21*, released in 2015, describes the educational policy for compulsory schools. It sets the educational goals at all school levels and informs all stakeholders about the competencies to be achieved in compulsory education. *Lehrplan 21* divides the eleven years of compulsory schooling into three cycles: (i) kindergarten, 1st and 2nd grade, (ii) 3rd through 6th grade, and (iii) 7th through 9th grade. The curriculum, furthermore, breaks the subjects down into competence areas, which focus on skills/abilities (e.g., listening, reading, speaking, writing in the languages) as well as thematic areas (e.g., “numbers and variables” in mathematics). Within these areas, competencies are defined in the form of typical “I can” statements, pointing to the abilities, which students are intended achieve at the end of each of the three cycles. The set of competencies in each of the subjects are ordered by seven competence levels, which are summarizing descriptions of the abilities and competencies the students hold. There is a strong relation to developmental theories (cf. Siegler et al., 2014): the levels are characterized by an increase in factual, conceptual, and procedural knowledge, by an increase in perceptive demands (e.g., speech comprehension), by increasingly complex application scenarios as well as the degree of self-regulation and independence that need to be applied. Related to *Bloom’s Taxonomy* of cognitive development (Anderson, 2013), a higher level encompasses the abilities and competencies of a lower level.

In a number of Swiss cantons (e.g., St. Gallen or Zürich), online-based test and training platforms are deployed; *Lernlupe* (www.lernlupe.ch) for 3rd – 6th grade and *Lernpass plus* (www.lernpassplus.ch) for 7th – 9th grade. These platforms provide individual training facilities and standardized online tests along the competencies and levels defined by *Lehrplan 21*. The feedback for students is formative and competence-oriented in nature. For example, students receive a verbal description of their abilities and their current competence level. The results of the standardized tests provide clear indications of strengths and existing competence gaps and they are utilized by the teachers to plan

an individual support of students. The tests, moreover, provide practical indications for career planning. For example, the *Jobskills* platform (www.jobskills.ch) enables a comparison of various job specifications with students' individual competence profiles.

Lernlupe and Lernpass plus feature IRT-based computer adaptive testing functions (cf. van der Linden, 2016). Adaptive testing allows optimizing the assessment quality within minimal testing times. The test items are selected on an individual basis so that the item difficulty matches the estimated ability of the student as exact as possible. The prerequisite for adaptive testing is extensive standardization studies to identify the item characteristics (e.g., item difficulties) based on a representative sample. In the cantons St. Gallen and Zürich, such large-scale studies (5000 students per age group) have been carried out in the past years. Based on these results, a metric between 200, which corresponds to the lowest level of difficulty or ability, and 800, which is the highest value, is established. The mean of this scale is 500. The scale is a "historic" IRT scale and used in a variety of standardized academic achievements tests, for example the GMAT (www.mba.com/exams/gmat).

The test and training platforms Lernlupe and Lernpass plus are used on a frequent basis by cantonal schools and accordingly rich is the basis of available data. The main user groups of the data are students, on the one hand, and teachers and parents on the other hand. The feedback formats generated by the systems are used, for instance, to inform parents about the learning progress of their children. The feedback is designed cautiously and restricted to the individual score (on the 200-800 scale) in relation to the achieved competence levels, accompanied by verbal categorical descriptions of strengths and achievements. An increasingly important aspect of data visualization refers to the identification of performance indicators and performance comparisons for administration and management on a local school level but also on a regional, political level. In comparison to the "traditional" methods of gathering data about the performance of schooling (e.g., the OECD PISA studies), the data from the aforementioned (and other) platforms are more up-to-date, rather longitudinal in nature, and more detailed. This increases the utility of the data significantly - and also the interest of local, regional, and national authorities.

2 VISUALIZING GENERAL PERFORMANCE DATA

To inform stakeholders about general learning performance, for example, of entire cohorts, a visualization type is necessary that includes all relevant information and that particularly offers a comparison of local data (the data of a specific school) with regional data and the data of the standardization studies. Relevant variables are cohort/age group, subject and competence areas, gender, as well as the temporal progression over years. The dependent variables are student performance on the 200-800 scale and the achieved competence levels. The challenge is to develop a form of visualization that conveys the meaning of these complex data in a very simple way.

For Lernlupe and Lernpass plus we designed a set of visualization formats including density diagrams and pie charts. A group of experts and users chose a scatter plot approach as the most intuitive form of visualization (cf. Figure 1). The main advantage of the scatter plot was seen in the fact that each student can be represented as an individual point. This was considered being highly intuitive to understand for a broad variety of users with a broad range of data literacy levels (i.e., children, parents, teachers, school leaders, administrative leaders, politicians).

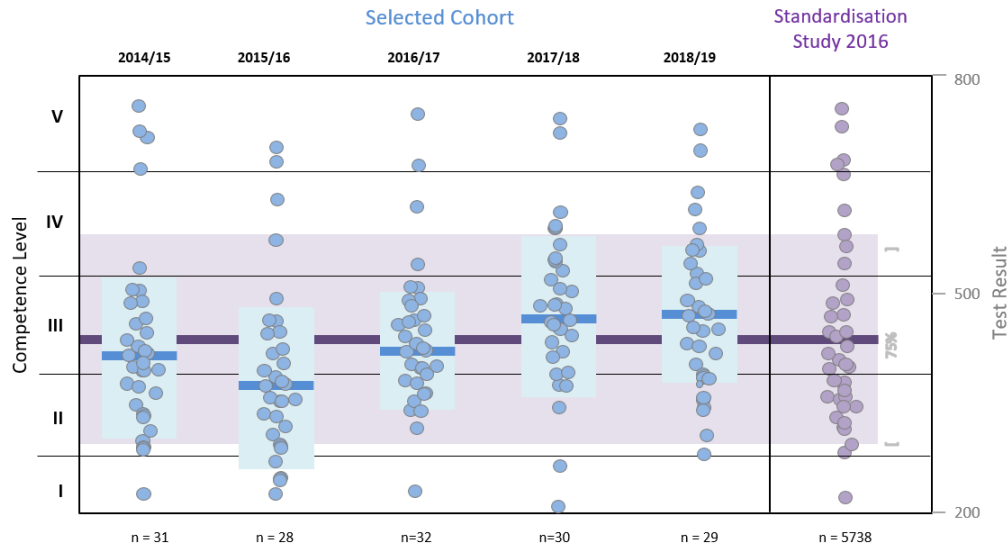


Figure 1: Visualizing learning data with a scatter plot. In this example, the performance of boys and girls is displayed on a yearly basis in comparison to the results of the standardization sample; the performance is shown as competence levels (scale on the left) and test scores (scale on the right).

3 JITTER AND SAMPLING TECHNIQUES

Scatterplots are amongst the most effective forms of understanding bivariate relationships, since they nicely display the relationship between a variable x and a variable y . However, typical scatter plots are only effective for two continuous variables. If one variable is discrete, other techniques for visualizing the data may be more appropriate. In our case, we have the continuous performance variables, however, only few classes (e.g., gender or years; cf. Figure 1). Also, we face the challenge that comparably small groups of students (e.g., 20 children of a class) are supposed to be compared to huge groups (e.g., 5000 students of the standardization studies). When displaying such a large amount of data points, the scatter plot gets illegibly crowded and confusing. Therefore, the visualization would lose its major strength. To overcome these issues, we developed a visualization approach combining sampling and jittering techniques to improve comprehensibility.

Jittering refers to adding random noise to data in order to prevent data points being over plotted by others. This over plotting specifically occurs when a continuous measurement (such as the standardized competence value) is rounded or aggregated. In large data sets, such as the standardization sample, over plotting is very likely and reduces the comprehensibility of the visualization substantially. Jittering can be done by adding small random changes to the actual values along the x or y -axes.

In our case, we have a small number of discrete classes (e.g., gender or the year), so we applied jittering along these classes (on the x -axis). A purely random jittering, however, may result in a low comprehensibility of the plot. We experienced that users tend to misinterpret the deviations from the center of a class. As a result, we developed an algorithmic jittering, which plots the data points with a minimum overlap to surrounding points. Given two points with exactly the same value, the

points overlap by a certain percentage. The more points exist with the same value, the higher is the percentage of overlapping (and therefore the smaller the visible area of the point) until a maximum overlap is reached. In a second step of the plotting approach, the same principle is applied to points with higher and lower values. By this means, the data points are not randomly jittered but iteratively placed within the outlines of the functional shape of the data (usually a bell curve). Technically, the basis for the algorithmic jitter function is *SinaPlot*, which is an enhanced jitter strip chart package in *R* where the width of the jitter is controlled by the density distribution of the data within each class (Sidiropoulos, et al. 2017, 2018). Figure 2 illustrates the approach.

Sampling refers to selecting only a (more or less representative yet small) sample from a pool of available data and to display this sample in the scatter plot. The technical challenge is to find the right sampling method for a given visualization purpose and a given target audience. In our example, sampling is done when the selected group of students exceeds 99 students and sampling is applied for comparisons with the standardization study.

One sampling method is to select random data points (i.e., students). Assuming that the original data follow a normal distribution, this results in a suitable representation of the original data. This approach, however, cannot guarantee that the most extreme students are represented, which is an important information. Also, this approach works well for very large data pools, such as the standardization group. For smaller pools, for example when a user selects the students of multiple classes, this method likely results in inadequate, most often a too uniform, sampling of observations. Our solution is a threshold-based sampling algorithm; for a specific number of observations (e.g., for 5 students with exactly the same value) only a single point is plotted. The algorithm also assures that for all individual data values, at least one point is maintained. When the number of points exceeds a maximum number (e.g., the aforementioned 99 points), the threshold is recursively raised. Moreover, the threshold follows a normal distribution, meaning that the threshold value is higher in middle areas than at the tail ends. This method allows displaying a distribution as close to the shape of the original one, with losing as little individual values as possible and without losing the most extreme values.

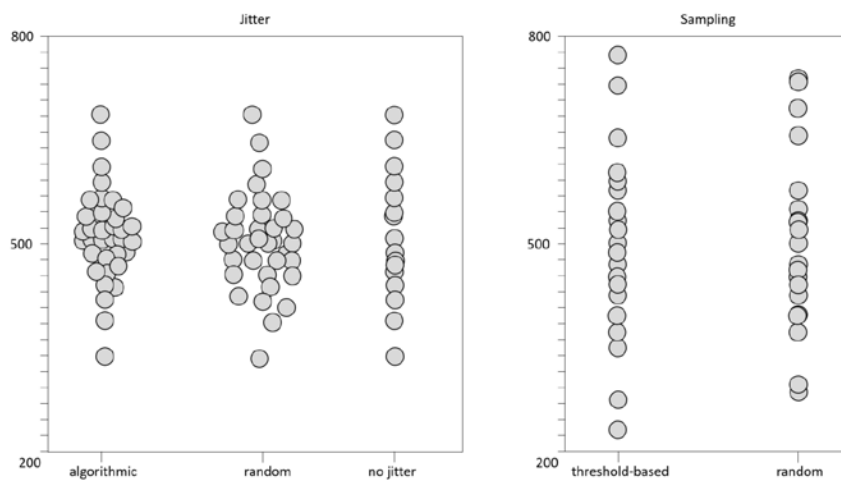


Figure 2: The left panel shows a comparison of the algorithmic jitter function and random jitter, opposed to the plot without jitter; the right panel shows a comparison of sampling methods

4 CONCLUSIONS

The approach of using jitter and data sampling turned out to be a suitable solution to make scatter plots more comprehensible and perhaps more applicable to wider scenarios. However, the algorithmic approach is arbitrary, in a way, because it bears a large degree of freedom. Therefore, finding the optimal settings for a specific set of data and use cases still is difficult. A critical factor, for example, is the screen resolution. The larger the screen, the more data points can and should be displayed. This, in turn, strongly influences the setup of the optimal plotting algorithm.

In user interviews with teachers and school leaders, the scatter plot was chosen as the most appropriate chart type to visualize the achievements of groups of students without losing the information about individuals. Even for large students groups (e.g., the standardization sample), the scatter plot was preferred over more conventional methods such as density functions, typical bell-shaped curves, or pie charts. The individual data points could easily be associated with “real” students, which was not the case with the rather abstract area below a curve, as an example. The downside of the scatter plot, particularly when larger amounts of data (i.e., > 50 points) are displayed on a typical computer screen, is a rapid decrease of comprehensibility and legibility, mainly due to an overlap of data points. I presented two techniques to maintain the strength of the scatter plot and reducing the issue of over plotting. This approach to display student data was the favorite among a group of potential users.

To explore different characteristics of the plotting algorithm for different scenarios, in further steps, we will conduct simulation studies to compare different configurations of the algorithm in terms of legibility and ease of comprehension.

REFERENCES

- Anderson, L. (2013). *A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives*. Cambridge, UK: Pearson Publishing.
- Sidiropoulos, N., Sohi S.H., Rapin N., & Bagger F.O. (2017). *An Enhanced Chart for Simple and Truthful Representation of Single Observations over Multiple Classes*. CRAN R package. Available online at <https://cran.r-project.org/web/packages/sinaplot/vignettes/SinaPlot.html>
- Sidiropoulos, N., Hadi Sohi, S., Pedersen, T.L., Porse, P.T., Winther, O., & Rapin, N. (2018). SinaPlot: An Enhanced Chart for Simple and Truthful Representation of Single Observations Over Multiple Classes. *Journal of Computational and Graphical Statistics*, 27(3), 673-676.
- Siegler, R. S., DeLoache, J. S., Eisenberg, N., & Saffran, J. (2014). *How Children Develop*, 4th edition. New York: Worth.
- van der Linden, W. (Ed.) (2016). *Handbook of Item Response Theory (Volume 1)*, 1st Edition. London, UK: Chapman and Hall/CRC;

Less (context) is more? Evaluation of a positioning test feedback dashboard for aspiring students.

Nicolas Hoppenbrouwers
Faculty of Engineering Science, KU Leuven

Tom Broos
Dept. of Computer Science, KU Leuven
tom.broos@kuleuven.be

Tinne De Laet
Tutorial Services, Faculty of Engineering Science, KU Leuven
tinne.delaet@kuleuven.be

ABSTRACT: Aspiring engineering students profit from feedback regarding how their mathematical skills compare to the requirements and expectations of an engineering bachelor program. The positioning test is a non-binding test used in Flanders, Belgium assessing the mathematical skills of aspiring students. This paper elaborates on the research on and development of a learning analytics dashboard (LAD) that provides feedback on a participants' obtained results. Its objective is to provide actionable insights and to raise awareness and reflection about the participants' strengths and weaknesses, and subsequently their choice of study. To reach the final dashboard, the design went through six iterations, 662 students were surveyed, and 60 persons were thoroughly interviewed, including study advisors, students, and visualization experts. The final dashboard was evaluated using the EFLA, SUS, and a custom-made questionnaire, and a framework of factual, interpretative, and reflective insights. The results show that the developed dashboard is a considerable improvement over a comparable state-of-the-art dashboard. Furthermore, results show that a more visual representation, confined to only the most essential information, provides a better overview, leads to more and deeper insights while displaying less information and context, and has better usability and attractiveness scores than a more textual version.

Keywords: learning analytics, information visualization, student dashboard, positioning test, learning technologies

1 INTRODUCTION

The first bachelor year is often cited as the most essential to future academic success [1, 2, 11, 28]. A wide range of research focuses on identifying predictors of academic success in the first bachelor year, before students enroll in university programs, as this would shed light on the skills and knowledge students need to be successful. Apart from the obtained grade-point average in secondary education [3, 29], literature often describes mathematical ability as the most significant predictor of persistence and attainment in STEM fields [18, 20, 22, 23]. Starting mathematical competences is identified as one of the primary factors determining whether a student will continue studying in a STEM field, and certainly for engineering [21, 27]. Once the relevant skills are identified, learning analytics dashboards (LAD) can be developed to provide aspiring students with feedback, hereby supporting them in the

transition from secondary to higher education (HE). LADs are an effective and commonly used tool in learning analytics (LA) to visualize information [5, 7, 14, 15, 26]. Just like the general objective of information visualization, they can be used to represent large and complex quantities of data in a simple way [15, 19]. Few [16] defines a dashboard as ‘a visual display of the most important information needed to achieve one or more objectives; consolidated and arranged on a single screen so the information can be monitored at a glance’. Unlike most other countries, students in Flanders do not have to complete any formal application procedure or test in order to enroll in a university program. Furthermore, the tuition fee of EUR 922 per year is relatively low compared to other nations. Consequently, students are free in their choice of program, resulting in a large degree of heterogeneity in the first bachelor year regarding knowledge, skills, and educational background. This results in a drop-out of 40% in STEM fields. Since 2011, the Flemish universities offering engineering bachelor programs have joined efforts for organizing the ‘positioning test’, a non-obligatory and non-binding diagnostic test for the candidate students’ ability to solve math problems [31]. The focus on mathematics is not surprising considering the importance of mathematical ability as a predictor for student success in STEM [18, 20, 22, 23]. The positioning test typically contains 30 multiple choice questions and is organized in the summer between the end of secondary education and the start of higher education.

This paper presents the research that aimed at developing a LAD that provides aspiring engineering students with feedback on their mathematical problem-solving skills, based on their results on the positioning test. The developed LAD aims at visually triggering insights in the obtained results. More specifically, the LAD should provide actionable insights, making students more aware of their strengths and weaknesses, and allowing students to reflect on their study choice. The objective of the LAD is similar to that of the positioning test itself, in that it tries to encourage and motivate students that do well on the positioning test (score > 14/20) to consider engineering as a viable and interesting study option, participants who obtain a low score (score < 8/20) to reflect on their study choice, and support the middle group to take remedial actions (e.g. a summer course) to improve their mathematical abilities in order to successfully attain an engineering degree. To achieve these objectives, the research ran through all phases of a user-centered design process, including a preliminary data-analysis, a large survey of 622 end users, pilot testing, and 55 in-depth interviews. Different evaluation metrics were used to assess the developed dashboard: EFLA [24, 25], SUS [4], and a custom-made questionnaire, and the framework of factual, interpretative, and reflective insights [10]. Finally, this paper compares the developed dashboard with an existing feedback dashboard [6] for the positioning test.

2 RELATED WORK

The literature describes several guidelines for developing effective LADs. For example, Few [16] describes thirteen commonly made mistakes when developing dashboards. Together with the general graphical integrity and design aesthetic principles defined by Tufte and Graves-Morris [30], they serve as the basis for the development of the dashboard. The most commonly used visualization types in LADs are bar charts, line graphs, tables, pie chart, scatterplot, simple text, world clouds and traffic lights. De Laet [12] however warns not to use traffic lights, and mentions how wording is essential in LA applications. Predictive LA applications have uncertainty and it is important this uncertainty is also displayed [12]. LADs should avoid speaking too much in terms of “chances of failure” and “success”

[12]. Two additional relevant guidelines are defined by Charleer et al. [8]. They recommend that LADs should be sufficiently aggregated or abstract as an uncluttered representation incites more detailed explorations of the LA data. Secondly, they recommend that LADs should provide functions that increase the level of detail in the data [8].

The LAD of this paper focuses on the transition from one education system to the other (secondary to HE), while most examples in the literature are more concerned with monitoring study progress during an educational program, either for a student or a tutor. Several LADs were used as an inspiration for the LAD of this paper, such as the OLI dashboard [13], the Course Signals dashboard [1], the Student Activity Meter (SAM) [17], and the LISSA-dashboard [9]. The most related dashboard is that state-of-the-art dashboard by Broos et al. [6], which also aims at providing feedback after the positioning test. This LAD referred further on to as the “reference dashboard” provides, beside feedback on the mathematical problem-solving skills of students, feedback on learning and studying skills, and the prior education of students [6]. The reference dashboard by Broos et al. contains elaborate textual explanations and feedback to contextualize the participants’ obtained result.

LADs can incorporate insights of other research while visualizing data. Vanderoost et al. [31] analyzed the predictive power of the positioning test for engineering studies in Flanders. More specifically, the research examines whether it is possible to “predict” first-year academic achievement using the results of the positioning test. More specifically, the goal is to identify three distinct groups of students: group A are students who perform well in their first bachelor year, achieving a study efficiency of over 80% after the first exam period in January; group C are with a study efficiency below 30%; group B are students with a SE between 30 and 80 %. Earlier research [31] showed that participants obtaining a high score on the positioning test ($>13/20$) more often obtain good study results (study efficiency (SE) $>80\%$) in the first semester (51%), while students with a low score on the positioning test ($<8/20$) more often do not enroll (35%), drop-out (6%), or have low academic achievement (SE $<30\%$) in the first semester (39%). Vanderoost et al. also showed how the study efficiency in the first semester of the first bachelor year strongly predicts if a student will complete the engineering bachelor and in which time frame (in 3, 4 or 5 (or more) years).

3 CONTEXT

The positioning test consists of approximately thirty multiple-choice questions assessing participants’ problem-solving skills. Formula scoring is used to calculate the overall result (on 20) based on each participant’s responses. Each question is assigned to one of five mathematical categories: (1) reasoning, (2) knowledge of concepts, (3) spatial visualization ability, (4) skills (calculating derivatives, solving systems of linear equations, combinatorics, geometry, etc.) (5) and modeling questions (problem solving questions in a physical context that need combination and modeling of different inputs). Additionally, each question is assigned to one of four difficulty levels. The difficulty level of a question is determined by the percentage of participants that correctly answered the question: the 25% best answered questions of the 30 questions have a difficulty level of 1, while the 25% worst answered questions have a difficulty level of 4.

End-users of the existing reference LAD are participants of the positioning test, consisting mainly of students that just completed secondary education. They receive access to their personalized LAD through a feedback email, typically three days after completing the test. Apart from these aspiring

engineering students, other stakeholders are also involved. The Tutorial Services of the Faculty of Engineering Science heavily participates in the development of the LAD.

They are represented by the head of the unit and two study advisors (SAs), who from their pedagogical experience and educational expertise give feedback on content and design. SAs are concerned with guiding and helping students with any questions they might have. They can also be considered end-users of the dashboard, as they use the LAD to start the conversations with participants that need more feedback and advice during a private meeting. LA researchers and visualization specialists, represented by three experts of an HCI research group, evaluate the quality of the design.

4 DESIGN

Design process. A user-centered design process was followed to develop the dashboard. The design passed six iterations before reaching its final state. Throughout the iterations, the design principles by Tufte and Graves-Morris [30], the commonly defined dashboard mistakes by Few [16] and a set of self-defined design requirements served as guidelines for the development of the dashboard. The self-defined design requirements are formal specifications of the general objective described in Section 1 identified based on interviews with the involved stakeholders. They consist of eight functional requirements and six non-functional requirements. An example of a functional requirement is: ‘the ability to compare your own result with the results of other participants’. An example of a non-functional requirement is: ‘a good balance between textual and visual elements’.

In total the dashboard was developed and improved in six iterations. Each iteration is characterized by a different objective, format, and evaluation method. The first iterations focused more on functional requirements, finding out expectations, and determining the right content. Later iterations focused more on non-functional requirements and correctly choosing and improving the visualizations. The final design was programmed using D3.js. Different methodologies were used for creation and evaluation of the dashboard, such as co-designing, rapid prototyping, guidelines and general principles, the EFLA and SUS questionnaire, formal presentations with feedback, and semi-structured as well as informal interviews, based on distinct protocols, for instance scenario-based with concurrent

The content of the dashboard has changed throughout the six iterations. We conducted semi-structured interviews with two study advisors of the Bachelor of Engineering Science at the Catholic University of Leuven (KU Leuven), informal interviews with the head of Tutorial Services of the faculty, and a questionnaire among 662 students. In the questionnaire, students scored 14 different content part suggestions on a 7-point Likert scale for relevance and usefulness to include in a feedback system after participation in the positioning test. Results show that students like to see their total score, a comparison to other participants, and the typical performance (in terms of SE) of first-year bachelor students that obtained a similar score on the positioning test the previous year. They also liked to see the aggregated score per mathematical category and the score and original assignment per question. Students were divided when it comes to displaying the typical performance (in terms of grades on the course) on each individual course of first-year bachelor students who obtained a similar score on the positioning test previous year. They also disagreed regarding the presence of a specific, personalized, pre-defined study choice advice, due to insufficient face validity. Confirmed by the results of a data-analysis, which showed a lack of predictive power for these features, we decided to remove them

from the dashboard. The conclusions of the interviews with the study advisors (SAs) are similar to those of the survey. Examples of features that were added throughout the iterative process are the aggregated score per degree of difficulty and a picture of the original question of the positioning test, as both study advisors and students reacted positively to these suggestions.

Final design. The final design of the dashboard exists in two variants, differing in one part. Fig. 1 displays the first variant and consists of five major parts. Each part has a tooltip presenting more detailed information, following the general guidelines proposed by Charleer et al. [8], described in Section 2. Fig. 3 shows the tooltips for part A and B of Fig. 1. Furthermore, a help icon on the top right corner of each graph contains more context and explanation, e.g. explaining the color of a graph. The five major parts of the LAD (Fig. 1) and its tooltips allow students to:

- (A) review their obtained total score and compare themselves to the other participants by showing the general score distribution of all participants;
- (B) review each question, its difficulty, its mathematical category, its original assignment, the answer they submitted and the correct answer;
- (C) review their aggregated obtained score per mathematical category or degree of difficulty, allowing them to find their strengths and weaknesses and see whether they score well/bad on certain categories or easier/harder questions, permitting them to discover whether they lack basic knowledge or only more advanced skills;
- (D) compare themselves to other participants for each mathematical category and degree of difficulty, by plotting the score distribution per topic ;
- (E) view the typical first-year academic achievement, in terms of SE, of students in earlier cohorts based on their total positioning test score, via a Sankey diagram.

The second variant, in part displayed in Fig. 2, differs only on the strengths & weaknesses section (part C and D in Fig. 1). It combines the information of these two parts in on large histogram, displaying the distribution of the total positioning test score of all participants, and five small histograms, displaying the score distribution per category or degree of difficulty. The objective of the two separate variants is to see which visualization leads to more insights and whether the click functionality of the first variant is intuitive.

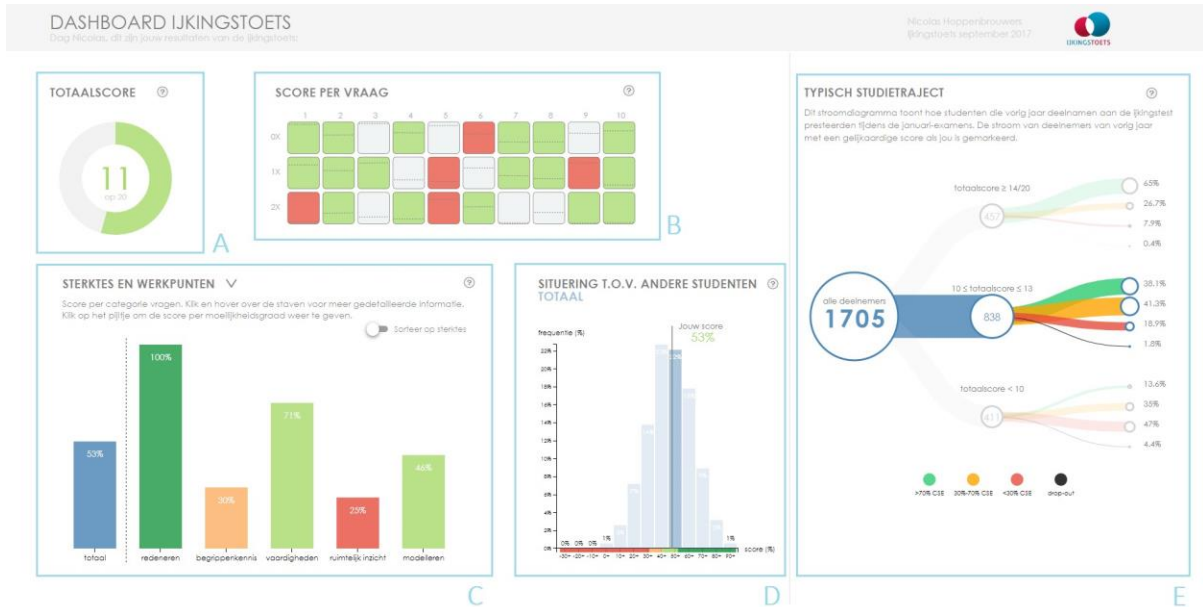


Figure 1: Feedback dashboard for future students after positioning test: first variant of final design. Corresponding to the displayed letters: (A) Donut chart showing the individually obtained total score on the position test (on 20). (B) Matrix showing the score and difficulty per question. Red indicates the student provided a wrong answer, grey for no answer and green for a correct answer. The percentage of correct responses by all test participants is indicated by the horizontal bar for each item. (C) Bar chart illustrating the participant’s strengths and weaknesses, by showing the score per mathematical category (reasoning, knowledge of concepts, spatial visualization ability, skills, and modeling) and per degree of difficulty. (D) Histogram showing performance of peers for each mathematical category and degree of difficulty. The student’s score is positioned using a vertical line. (E) Sankey diagram showing performance of previous students in the first bachelor year with a comparable total score on the positioning test. The diagram shows the outcomes in study efficiency (e.g. green arrows for students achieving 70% of study points, black arrows for students dropping out of the program) for three groups of positioning test outcomes (less than 10/20, from 10 to 13 and above 13).

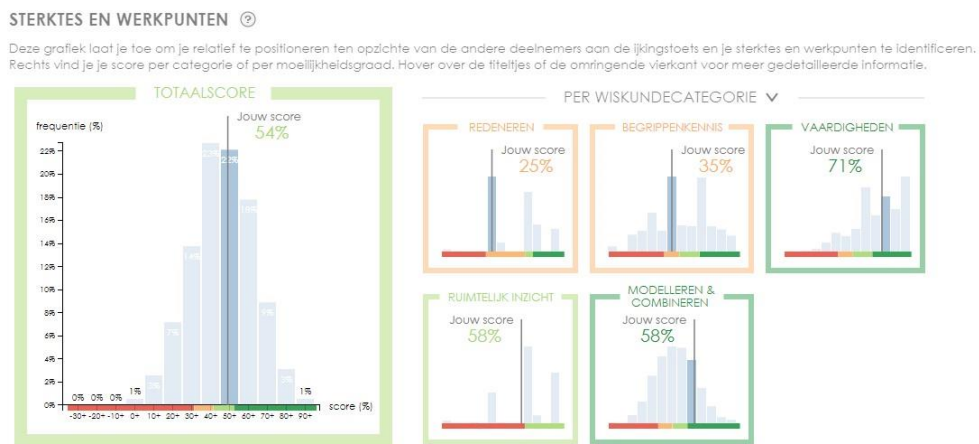


Figure 2: An alternative visualization for the score per category in the final dashboard, substituting part C and D of the dashboard.

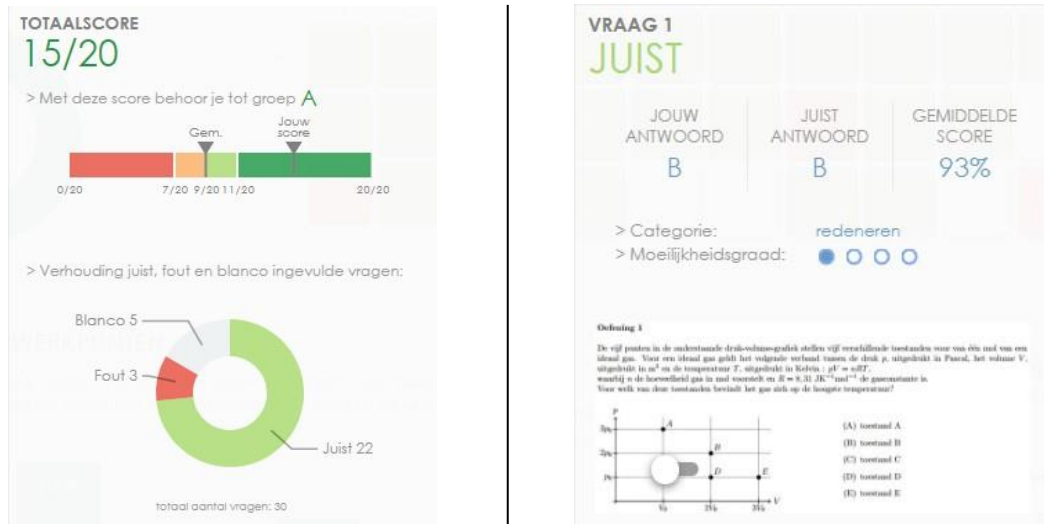


Figure 3: Two examples of tooltips in the final dashboard. The tooltip on the left is shown when clicking or moving the mouse over part A in Figure 1. It shows more detailed information about the performance of the student on the positioning test and compares the result to that of other participants. The tooltip on the right clicking or moving the mouse over any of the items in part B in Figure 1. It shows detailed information for each question on the test, including the given and correct answer, the difficulty level, and question category. At the bottom of the tooltip, the question is represented in the same way as in the positioning test (on paper) to aid visual recall.

5 EVALUATION

Both variants, described in Section 4, are evaluated and compared to the reference dashboard [6], described in Section 2. The latter is currently in use in the Flemish universities organizing the positioning test. It is text-heavy in comparison to the evaluated alternatives, which allows to assess the added value of inclusion of such elaborate textual guidance.

Evaluation of the two final variants of the dashboard and reference dashboard [6] is based on 48 in-depth interviews (16 per dashboard), each lasting between 40 minutes and 1 hour. Each interview consists of four stages. The first phase of the interview is scenario-based, using the concurrent think-aloud protocol. End-users have to imagine having participated in the positioning test and now getting their result. Three scenarios are possible. Either they get a score in which they belong to group A (total score of 15/20), either group B (12/20) or group C (6/20). Anonymized data is used from the dataset described in Section 4. Each test user says out loud the insights they obtain upon visualization of the dashboard. The framework by Claes et al. [10] is used to measure these insights. The framework defines three levels of insights: 1) factual insights: simple, objective statements or questions that are triggered by the dashboard, e.g. “I obtained a total score of 14/20.”; 2) interpretative insights: interpretation of the displayed data, relying on the participant’s knowledge and experiences, e.g. “I mainly score well on the easier questions.”; 3) reflective insights: subjective, emotional and personal connotations triggered by the dashboard, leading to further awareness and reflection, e.g. “I feel like I did not do well enough at this test, making me doubt about whether I should go for another study program.”. Each insight is categorized into one of these levels. The test user can also mention when something in the dashboard is unclear, but the monitor of the test does not intervene and only writes down all statements made by the test person.

In the second phase, the interview switches to a task-based interview with active intervention. The monitor gives the test persons tasks based on the information or insights they missed during the first phase and finds out why these parts and insights have been missed. This phase tries to examine whether the dashboard is intuitive and has any shortcomings.

In the third phase, the test person fills in the SUS, the EFLA and a custom-made questionnaire, which verifies whether design requirements have been met. The EFLA questionnaire has been translated to Dutch and adapted to reflect the topic of the dashboard, identical to the evaluation of the dashboard of Broos et al. [6]. The design requirements questionnaire test consisted of 21 statements, to which the user could “Strongly disagree” or “Strongly agree”, using a 5-point Likert scale.

Finally, in the fourth phase the test persons get to see the two other dashboards and can express their preference. This last phase was optional.

6 RESULTS

Based on the recorded insights during the interviews 13 types of factual, 11 of reflective, and 8 types of interpretative insights were identified. All types of insights occurred more often with the participants for the LAD developed in this research compared to the reference dashboard (Table 1).

Table 1: Subset of the 13 types of factual (F), 11 types of reflective (R), and 8 types of interpretative (I) insights identified during the interviews and the percentages of interviewees in which these insights were found for the reference dashboard (B) of [6] and the two variants described in this paper.

Description	insight	%B	%V1	%V2
(F1)	My total score on the positioning test was ...	100	100	100
(F2)	My total score placed me in group A/B/C ...	100	94	94
(F3)	I answered X questions correct/wrong/blank	75	100	100
(F4)	My answer to this question was correct/wrong/blank	81	100	100
(F5)	On average this question was replied well/badly	56	88	94
(I1)	My total score compared wrt other participants	100	100	100
(I2)	This question was difficult/easy	56	88	81
(I5)	I score especially well in easy/difficult questions	56	56	63
(R1)	Reflection on total score	100	100	100
(R2)	Reflection on comparison wrt peers	69	100	94
(R3)	I guessed/left blank too many questions	44	56	63
(R4)	Reflection on particular question	56	88	81
(R10)	Reflection on future academic achievement	69	88	94
(R11)	Reflection on study choice	75	100	94

Fig. 4 shows the total SUS and EFLA score and the score per EFLA-dimension. The first variant has an overall average SUS-score of 81, the second variant 76, both statistically significant ($p < 0.01$) higher than the score of 47 of the reference dashboard. A score of more than 68 is considered above average [4], implying that the developed LAD has a better usability design than the reference dashboard. The differences between the averages of the two variants of the final dashboard are not statistically significant ($p > 0.2$). The total EFLA-score of the first variant is 74 and of the second variant is 70. Only

the EFLA score of the first variant is statistically significantly higher than the one of the reference dashboard score of 59.

4

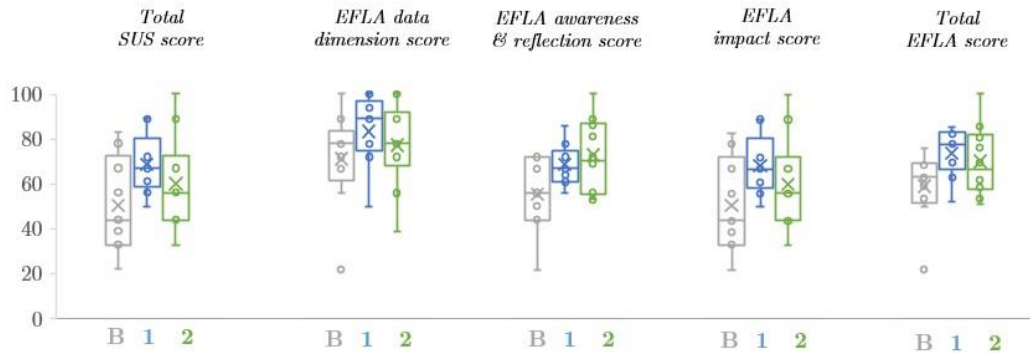


Figure 4: The total SUS and EFLA score and the score per EFLA-dimension: the data dimension (questions D1+D2), the awareness and reflection dimension (A1-A4) and the impact dimension (I1+I2). Gray boxplots ('B') denote the reference dashboard [6], blue box-plots ('1') denote the first variant of the final design of this paper and green ('2') the second variant.

The results of the design requirements questionnaire showed that each of the three dashboards successfully helps participants in understanding whether their current mathematical skills are matched with the expected mathematical skills and incites users of the LAD to awareness and reflection. Both variants, however, scored significantly better than the reference dashboard on the ability to use the dashboard independently, give a better overview of strengths and weaknesses, give a better detailed overview of the obtained result and allow participants to compare themselves more to the other participants. The users also indicated that these dashboards are better at displaying only factual, objective information, without giving interpretations or conclusions, but indicated that the dashboards can also be more confronting. Furthermore, they found that the two variants were more personalized, immediately gave an indication of the most important information, were better at showing only information that is relevant, were better at providing context, were more aesthetically pleasing, add less ambiguity and have a better balance between textual and visual elements, compared to the reference dashboard. For most design requirements, the differences between the two variants are not statistically significant.

7 DISCUSSION AND CONCLUSION

7.1 Implications for LAD design

This dashboard provides feedback to participants of the positioning test for the engineering program, inciting awareness and reflection about their strengths and weaknesses, and consequently their choice of study. The results of this LAD are interesting, as it focuses on the transition from secondary

school to higher education, while most LADs in the literature focus on monitoring students when they are already at university or college. Furthermore, a comparison has been made with the reference dashboard [6] that is currently used for feedback to the participants of the positioning test. The LADs developed in this research are more visual compared to the reference dashboard. Following thorough evaluation of the six iterations of the dashboard, the most important advantages of the more visual dashboards in this paper are that they have better usability scores, provide a better overview of the obtained results and a participant's strengths and weaknesses and visualize only relevant and objective information. A surprising result is that, while the visual dashboards contain less context and explanation, they still lead to more interpretative and reflective insights. Users declare that they think the layering of detail is better in the more visual dashboards. The main screen provides a good overview and immediately gives an indication of the essence, while the tooltips allow for more detailed information, consistent with the guidelines of Charleer et al. [8]. According to the tests, the reference dashboard of Broos et al.[6] has too much unnecessary information and text, which leads to users getting lost and not knowing what they should learn as take-away message. Some test persons also admit skipping parts of this dashboard because they "do not want to read so much text", causing them to miss out on important information.

The first most important general conclusion is that confining LADs to the most essential information, not displaying an overload of context and explanations, but using intuitive and simple visualizations, displaying less information, may lead to more awareness and reflections. An important part of LA applications is to make sure the end-users cannot get the incorrect interpretation, often leading to a lot of textual clarification. This research tries to convey to the designer that more text not necessarily means better insights, but well-designed and intuitive visualizations do.

Secondly, many test users mention how the dashboards of this paper are aesthetically pleasing and "fun to play with". Animations direct the user's attention to the most important information but are also specifically included to make the dashboard more aesthetically pleasing and show that the data is dynamic and interactive. While this result seems only of minor importance, it should not be underestimated. Several users mention how the aesthetics make them want to play more with the dashboard and spend more time with the dashboard. This eventually leads to more insights, which is essentially the final goal of this LAD. A lot of LADs do not spend enough time on the aesthetics of the dashboard, underestimating the effect this has on the effectiveness of the dashboard.

Finally, another objective was to see which of the two variants is more effective. The differences in the results are however not statistically different. Most users prefer the first variant, as it seems less cluttered at first sight, but end-users often miss some of the functionality in this variant. Further iterations should combine the best elements of both visualizations.

7.2 Future work and limitations of the study

The more visual dashboards however also have several disadvantages and pose new challenges. As all information is displayed on a single screen, some users observe the dashboard in an unstructured way, sometimes leading to less interpretative or reflective insights and confusion. Most participants observed the dashboard in a structured manner, but further research could examine whether a different arrangement of the various graphs could resolve this issue, keeping the visual character of the dashboard. Suggestions are a more sequential ordering of the graphs, similar to a grade report in

high school, or to use a guided tour to force the correct logical flow. Secondly, extra care is needed for the placement and highlighting of text. Because the visual dashboard looks more intuitive, users are less inclined to read any text at all, acknowledged by several test persons. While the graphs are mostly clear by themselves and lead to more interpretative and reflective insights, this a real concern for the development of a dashboard. Further research should examine how to highlight text to force the user's attention to the surrounding text, even if they already understand the graph.

This study presents both qualitative and quantitative results of thorough four-stage evaluations with test users. It must be noted that the evaluation of the LADs happened with more experienced students asked to imagine being in the randomly assigned scenario of a student in transition from secondary to higher education. Test users completed the SUS, EFLA and custom questionnaires after an in-depth and a task-based interview (see Section 6). This may contribute to the explanation of inter-study differences between results reported previously [6] for the reference LAD (overall EFLA score of 72) and those reported in this paper (overall EFLA score of 59). In the former study, the actual target group of the reference LAD was surveyed using an on-screen questionnaire available within the dashboard itself. Further work is necessary to assess if, once accounted for methodological influence, outcome differences indicate that experienced students have different needs and preferences for LADs than newcomers.

REFERENCES

- [1] Arnold, K. E., & Pistilli, M. D. (2012, April). Course signals at Purdue: Using learning analytics to increase student success. In *Proceedings of the 2nd international conference on learning analytics and knowledge* (pp. 267-270). ACM.
- [2] Besterfield-Sacre, M., Atman, C. J., & Shuman, L. J. (1997). Characteristics of freshman engineering students: Models for determining student attrition in engineering. *Journal of Engineering Education*, 86(2), 139-149.
- [3] Bridgeman, B., McCamley-Jenkins, L., & Ervin, N. (2000). Predictions of freshman grade-point average from the revised and recentered SAT® I: Reasoning Test. *ETS Research Report Series*, 2000(1), i-16.
- [4] Brooke, J. (1996). SUS-A quick and dirty usability scale. *Usability evaluation in industry*, 189(194), 4-7.
- [5] Broos, T., Peeters, L., Verbert, K., Van Soom, C., Langie, G., & De Laet, T. (2017, July). Dashboard for actionable feedback on learning skills: Scalability and usefulness. In *International Conference on Learning and Collaboration Technologies* (pp. 229-241). Springer, Cham.
- [6] Broos, T., Verbert, K., Langie, G., Van Soom, C., & De Laet, T. (2018, March). Multi-institutional positioning test feedback dashboard for aspiring students: lessons learnt from a case study in flanders. In *Proceedings of the 8th International Conference on Learning Analytics and Knowledge* (pp. 51-55). ACM.
- [7] Broos, T., Verbert, K., Langie, G., Van Soom, C., & De Laet, T. (2017). Small data as a conversation starter for learning analytics: Exam results dashboard for first-year students in higher education. *Journal of Research in Innovative Teaching & Learning*, 10(2), 94-106.
- [8] Charleer, S., Klerkx, J., Duval, E., De Laet, T., & Verbert, K. (2016, September). Creating effective learning analytics dashboards: Lessons learnt. In *European Conference on Technology Enhanced Learning* (pp. 42-56). Springer, Cham.
- [9] Charleer, S., Moere, A. V., Klerkx, J., Verbert, K., & De Laet, T. (2018). Learning analytics dashboards to support adviser-student dialogue. *IEEE Transactions on Learning Technologies*, 11(3), 389-399.
- [10] Claes, S., Wouters, N., Slegers, K., & Vande Moere, A. (2015, April). Controlling in-the-wild evaluation studies of public displays. In *Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems* (pp. 81-84). ACM.
- [11] Crisp, G., Nora, A., & Taggart, A. (2009). Student characteristics, pre-college, college, and environmental factors as predictors of majoring in and earning a STEM degree: An analysis of students attending a Hispanic serving institution.

- [12] De Laet, T (2018). The (non)sense of “chances of success” and predictive models. <http://blog.associatie.kuleuven.be/tinnedelaet/the-nonsense-of-chances-of-success-and-predictive-models/>. Accessed 4 April 2018.
- [13] Dollár, A., & Steif, P. S. (2012). Web-based statics course with learning dashboard for instructors. *Proceedings of computers and advanced technology in education (CATE 2012), Napoli, Italy*.
- [14] Duval, E. (2011, February). Attention please!: learning analytics for visualization and recommendation. In *Proceedings of the 1st international conference on learning analytics and knowledge* (pp. 9-17). ACM.
- [15] Elias, T. (2011). Learning analytics. *Learning*, 1-22.
- [16] Few, S. (2006). Information dashboard design.
- [17] Govaerts, S., Verbert, K., Duval, E., & Pardo, A. (2012, May). The student activity meter for awareness and self-reflection. In *CHI'12 Extended Abstracts on Human Factors in Computing Systems* (pp. 869-884). ACM.
- [18] Green, A., & Sanderson, D. (2018). The roots of STEM achievement: An analysis of persistence and attainment in STEM majors. *The American Economist*, 63(1), 79-93.
- [19] Khalil, M., & Ebner, M. (2016). What is learning analytics about? A survey of different methods used in 2013-2015. *arXiv preprint arXiv:1606.02878*.
- [20] Kokkelenberg, E. C., & Sinha, E. (2010). Who succeeds in STEM studies? An analysis of Binghamton University undergraduate students. *Economics of Education Review*, 29(6), 935-946.
- [21] Leuwerke, W. C., Robbins, S., Sawyer, R., & Hovland, M. (2004). Predicting engineering major status from mathematics achievement and interest congruence. *Journal of Career Assessment*, 12(2), 135-149.
- [22] Moses, L., Hall, C., Wuensch, K., De Urquidi, K., Kauffmann, P., Swart, W., ... & Dixon, G. (2011). Are math readiness and personality predictive of first-year retention in engineering?. *The Journal of psychology*, 145(3), 229-245.
- [23] Pinxten, M., Van Soom, C., Peeters, C., De Laet, T., & Langie, G. (2017). At-risk at the gate: prediction of study success of first-year science and engineering students in an open-admission university in Flanders—any incremental validity of study strategies?. *European Journal of Psychology of Education*, 1-22.
- [24] Scheffel, M (2018). Evaluation Framework for LA (EFLA). <http://www.laceproject.eu/evaluation-framework-for-la> , Accessed 1 March 2018.
- [25] Scheffel, M., Drachsler, H., & Specht, M. (2015, March). Developing an evaluation framework of quality indicators for learning analytics. In *Proceedings of the Fifth International Conference on Learning Analytics And Knowledge* (pp. 16-20). ACM.
- [26] Schwendimann, B. A., Rodriguez-Triana, M. J., Vozniuk, A., Prieto, L. P., Boroujeni, M. S., Holzer, A., ... & Dillenbourg, P. (2017). Perceiving learning at a glance: A systematic literature review of learning dashboard research. *IEEE Transactions on Learning Technologies*, 10(1), 30-41.
- [27] Tinto, V. (2005). *College student retention: Formula for student success*. Greenwood Publishing Group.
- [28] Solberg Nes, L., Evans, D. R., & Segerstrom, S. C. (2009). Optimism and College Retention: Mediation by Motivation, Performance, and Adjustment 1. *Journal of Applied Social Psychology*, 39(8), 1887-1912.
- [39] Stumpf, H., & Stanley, J. C. (2002). Group data on high school grade point averages and scores on academic aptitude tests as predictors of institutional graduation rates. *Educational and Psychological Measurement*, 62(6), 1042-1052.
- [30] Edward Tufte and P Graves-Morris. 1983. *The visual display of quantitative information*. Vol. 2. Cheshire, CT Graphics press.
- [31] Vanderost, J., Van Soom, C., Langie, G., Van den Bossche, J., Callens, R., Vandewalle, J., & De Laet, T. (2015, June). Engineering and science positioning tests in Flanders: powerful predictors for study success?. In *Proceedings of the 43rd Annual SEFI Conference* (pp. 1-8).

AutoTutor Tutorial: Conversational Intelligent Systems and Learning Analytics

Bor-Chen Kuo^{1*}, Chen-Huei Liao¹, Kai-Chih Pai¹, Chia-Hua Lin¹, Xiangen Hu^{2,3}, Zhiqiang Cai², Art Graesser².

¹National Taichung University of Education, Taiwan

²University of Memphis, USA

³Central China Normal University, China

* kbc@mail.ntcu.edu.tw

ABSTRACT: Conversational Intelligent tutoring system is a class of Adaptive Instructional Systems that are among the most studied and efficiently implemented in the last 20 years. This tutorial will introduce the most successful example C-ITS called AutoTutor and focuses on the authoring of AutoTutor lessons and Data analysis process of Tutoring data. Authoring of AutoTutor lessons include a) implementing discourse strategies in AutoTutor dialogues and trialogues, b) creating conversation elements (such as media elements); c) conversation rules, and d) using existing well-made authoring templates. Data analysis process of tutoring data include applying learning analytics methods, such as Bayesian Knowledge Tracing (BKT), Additive Factors Model (AFM), ...etc., to leverage the sequences of observations from student-ITS interaction log files to continually update the estimate of student latent knowledge.

Keywords: AutoTutor, Student Models, Learning Analytics

1 TUTORIAL BACKGROUND

Institute of Electrical and Electronics Engineers (IEEE) recently approved a standard committee ([P2247.1 - Standard for the Classification of Adaptive Instructional Systems](#)). This is a significant milestone for advanced personalized learning, which is identified by the National Academy of Engineering one of the grand challenges of the 21st century (<http://www.engineeringchallenges.org/9127.aspx>). Conversational Intelligent tutoring systems (C-ITS) is a class of AIS that are among the most studied and efficiently implemented in the last 20 years. This tutorial will bring you the most successful example C-ITS called AutoTutor (Graesser, Hu, & Person, 2001; Graesser et al., 2004; Nye, Graesser, & Hu, 2014; Nye, Graesser, Hu, & Cai, 2014; Person et al., 2000). AutoTutor holds conversations with the human in natural language. The authors of the proposed tutorial are among those who have development multiple versions of AutoTutor that teaches Critical Thinking (Wallace et al., 2009), Computer Literacy (Person, 2003), Physics (Graesser et al., 2003), Reading (Graesser et al., 2016), Electronics (Morgan et al., 2018), Chinese reading and mathematics learning (Liao, Kuo, & Pai, 2012).

AutoTutor applications are built with the guidance of human learning principles (A. C. Graesser, Halpern, & Hakel, 2008), such as Deep Questioning, to help students learn by holding deep reasoning conversations (Arthur C. Graesser & Person, 1994). AutoTutor converses with learners follow the Expectation-Misconception Tailored (EMT) dialog (Arthur C. Graesser et al., 2004). An AutoTutor conversation often starts with a main question about a certain topic. The goal of the conversation is to help students' construct an acceptable answer (expected answers) to the main question. Instead of telling the students the answers, AutoTutor asks a sequence of questions (hints, prompts) that target specific concepts involved in the ideal answer to the main question. AutoTutor systems respond to students' natural language input, as well as other interactions, such as making a choice, arranging some objects in the learning environment, etc.

This tutorial focuses on the authoring of AutoTutor lessons and Data analysis process of Tutoring data:

1. Authoring of AutoTutor lessons include a) implementing discourse strategies in AutoTutor dialogues and triologues, b) creating conversation elements (such as media elements); c) conversation rules, and d) using existing well-made authoring templates.
2. Data analysis process of tutoring data include applying learning analytics methods, such as Bayesian Knowledge Tracing (BKT), Additive Factors Model (AFM), ...etc., to leverage the sequences of observations from student-ITS interaction log files to continually update the estimate of student latent knowledge.

2 ORGANIZATIONAL DETAILS

This event is a full-day tutorial. A Moodle website will be set up to continuously add more details and materials for participants. Tutorial will be announced on AutoTutor website (autotutor.org). Participants need to bring laptops. An example AutoTutor lesson will be provided to participants. Participants will create one's own AutoTutor lesson by modifying the example lesson. The proposed agenda is presented below in Table 1.

Table 1: Proposed agenda

Time	Session	Content
9:00-9:15	Introduction to AutoTutor	Introduction – Introduction of presenters and participants
9:15-10:30		Overview and Demo of AutoTutor Systems
10:30-11:00	Coffee Break	
11:00-12:30	AutoTutor Script Authoring Tool	A step by step guidance to creating an AutoTutor lesson
12:30-14:00	Lunch Break	
14:00-15:30	Student Models and Learning Analytics	Student Models in AutoTutor
16:00-17:30	Learning Analytics for AutoTutor	AutoTutor log data analysis by using Bayesian Knowledge Tracing (BKT) and Additive Factors Model (AFM)

3 TUTORIAL OBJECTIVES OR INTENDED OUTCOMES

The objectives of the tutorial include but not limited to 1) Understand the theoretical foundations, enabling technologies, and practical applications of C-ITS through hands-on and worked-out examples of AutoTutor. 2) Familiar with simple, common, and advanced data analysis methods that apply to the analysis of AutoTutor Data.

The intended outcomes of the tutorial include 1) All attendees will be able to create a complete C-ITS module. 2) All attendees will understand the data structure of the interaction between C-ITS and learners, 3) All attendees will be able to analyze data using the data analytical methods introduced.

REFERENCES

- Graesser, A. C., Cai, Z., Baer, W. O., Olney, A. M., Hu, X., Reed, M., & Greenberg, D. (2016). Reading comprehension lessons in AutoTutor for the Center for the Study of Adult Literacy. *Adaptive Educational Technologies for Literacy Instruction*, 288–293.
- Graesser, A. C., Halpern, D. F., & Hakel, M. (2008). *25 principles of learning*. Task Force on Lifelong Learning at Work and at Home Washington, DC.
- Graesser, A. C., Hu, X., & Person, N. K. (2001). Teaching with the help of talking heads. *Proceedings of the 2001 IEEE International Conference on Advanced Learning Technologies*, 460-461.
- Graesser, A. C., Jackson, G. T., Matthews, E. C., Mitchell, H. H., Olney, A., Ventura, M., et al. (2003). Why/AutoTutor: A test of learning gains from a physics tutor with natural language dialog. In R. Alterman & D. Hirsh (Eds.), *Proceedings of the 25th Annual Conference of the Cognitive Science Society* (pp. 1-5). Boston: Cognitive Science Society.
- Graesser, A. C., Lu, S., Jackson, G. T., Mitchell, H. H., Ventura, M., Olney, A. M., et al. (2004). AutoTutor: A tutor with dialogue in natural language. *Behavior Research Methods, Instruments, & Computers*, 36, 180-193.
- Graesser, A. C., & Person, N. K. (1994). Question asking during tutoring. *American Educational Research Journal*, 31, 104–137.
- Liao, C.-H., Kuo, B.-C., & Pai, K.-C. (2012). Effectiveness of Automated Chinese Sentence Scoring with Latent Semantic Analysis. *Turkish Online Journal of Educational Technology-TOJET*, 11(2), 80–87.
- Morgan, B., Hampton, A. J., Cai, Z., Tackett, A., Wang, L., Hu, X., & Graesser, A. C. (2018). Electronixtutor Integrates Multiple Learning Resources to Teach Electronics on the Web. *In Proceedings of the Fifth Annual ACM Conference on Learning at Scale* (pp. 33:1–33:2). New York, NY, USA: ACM.
- Nye, B.D., Graesser, A.C., & Hu, X. (2014). AutoTutor and family: A review of 17 years of natural language tutoring. *International Journal of Artificial Intelligence in Education*, 24(4), 427-469.
- Nye, BD, Graesser, AC, Hu, X. (2014b). AutoTutor in the cloud: a service-oriented paradigm for an interoperable natural-language ITS. *Journal of Advanced Distributed Learning Technology*, 2(6), 35–48.
- Person, N. K. (2003). AutoTutor improves deep learning of computer literacy: Is it the dialog or the talking head. *Artificial Intelligence in Education: Shaping the Future of Learning through Intelligent Technologies*, 97, 47.

- Person, N. K., Craig, S., Price, P., Hu, X., Gholson, B., Graesser, A. C., & Tutoring Research Group. (2000). Incorporating human-like conversational behaviors in AutoTutor. *Proceedings of the Agents 2000 Conference*, 85-92.
- Wallace, P., Graesser, A. C., Millis, K., Halpern, D., Cai, Z., Britt, M. A., Magliano, J., & Wiemer, K. (2009). Operation ARIES!: A computerized game for teaching scientific inquiry. In V. Dimitrova, R. Mizoguchi, B. Du Boulay, & A. C. Graesser (Eds.), *Proceedings of 14th International Conference on Artificial Intelligence in Education. Building Learning Systems that Care: From Knowledge Representation to Affective Modelling* (pp. 602-604). Amsterdam: IOS Press.