

Immediate Elaborated Feedback Personalization in Online Assessment

Ekaterina Vasilyeva^{1,2}, Paul De Bra², and Mykola Pechenizkiy²

¹ Department of Computer Science and Information Systems, University of Jyväskylä,
P.O. Box 35, 40351 Jyväskylä, Finland

² Department of Mathematics and Computer Science, Eindhoven University of Technology,
P.O. Box 513, 5600MB Eindhoven, the Netherlands
{e.vasilyeva,m.pechenizkiy}@tue.nl, debra@win.tue.nl

Abstract. Providing a student with feedback that is timely, most suitable and useful for her personality and the performed task is a challenging problem of online assessment within Web-based Learning Systems (WBLSs). In our recent work we suggested a general approach of feedback adaptation in WBLS and through a series of experiments we demonstrated the possibilities of tailoring the feedback that is presented to a student as a result of her response to questions of an online test, taking into account the individual learning styles (LS), certitude in a response and correctness of this response. In this paper we present the result of the most recent experimental field study where we tested two feedback adaptation strategies in real student assessment settings (73 students had to answer 15 multiple-choice questions for passing the midterm exam). The first strategy is based on the correctness and certitude of the response, while the second strategy takes student LS into account as well. The analysis of assessment results and students' behaviour demonstrate that both strategies perform reasonably well, yet the analysis also provide some evidence that the second strategy does a better job.

Keywords: feedback authoring, feedback personalization, learning styles, online assessment, response certitude.

1 Introduction

Online assessment becomes an important component of modern education. Nowadays it is used not only in e-learning, but also within blended learning, as part of the learning process. Online assessment is utilized both for self-evaluation and for "real" exams and it tends to replace or complement traditional methods of evaluation of the student's performance.

Providing formative and summative feedback is especially crucial in online assessment as students need to be informed about the results of their (current and/or overall) performance. The existing great variety of the feedback functions and types that the system can actually support make the authoring and design of the feedback in e-learning rather complicated [13]. An important issue is that different types of feedback can have a different effect (positive or negative) on learning and interaction processes [3]. Badly designed feedback (and/or the lack of feedback) could distract

the student from learning; it could provoke the students to stop using the e-learning system or even to drop the course (even in blended learning).

Feedback adaptation and personalization [10] is aimed to provide a student with the feedback that is most suitable and useful for his/her personality, the performed task and environment. The development of the personalized feedback requires having the answers to at least the following questions: what can be personalized in the feedback; and to which user or performance characteristics feedback should be personalized. Some answers to these fundamental issues can be found in [13].

In this paper we present the results of the experimental study where we tested two immediate elaborated feedback (EF) adaptation strategies in the online assessment of students through multiple-choice quiz within the (slightly altered) *Moodle* WBLs. In the quiz, students had to select their confidence (certainty) level and were able to receive different (adaptively selected and recommended) kinds of immediate EF for the answered questions. Our first strategy is based on the analysis of response correctness and response certitude only, while the second strategy, besides the analysis of the response, takes student's LS into account as well.

The analysis of the assessment data demonstrates that both strategies perform reasonably well. The results of our analysis however favor the second strategy and thus advocate the benefits of taking into account LS for selecting and recommending the most appropriate type of EF during the online assessment.

2 Tailoring Feedback in Online Assessment in WBLs

Feedback may have different learning effects in WBLs; it can inform the student about the correctness of his responses, "fill the gaps" in the student's knowledge by presenting information the student appears not to know, and "patch the student's knowledge" by trying to overcome misconceptions the student may have [4, 5, 7].

The functions of the feedback imply the complexity of information that can be presented in immediate feedback: verification and EF [6]. Verification can be given in the form of knowledge of response (indication of whether the answer was received and accepted by the system), knowledge of results (KR) (correctness or incorrectness of the response), or knowledge-of-correct response (KCR) (presentation of the correct answers) feedback. With EF the system besides (or instead of) presenting the correct answer, provides also additional information – corresponding learning materials, explanations, examples, etc [9].

Different types of feedback can be differently effective (and can even be disturbing or annoying to the student thus having also negative influence) in learning and interaction [3]. E.g., an important issue in designing feedback is that it can draw attention away from the tasks, thereby increasing the time required to execute them.

Design of feedback assumes that the following questions can/must be answered: (1) when should the feedback be presented; (2) what functions should it fulfil; (3) what kind of information should it include; (4) for which students and in which situations would it be most effective? The variety of possible answers to these questions makes design of feedback rather complicated, especially in WBLs.

Our recent studies [8, 10, 11, 12] were aimed at demonstrating the feasibility and benefits of designing adaptive feedback (with respect to the characteristics of an individual student) in online multiple-choice tests.

Adaptive feedback is aimed at providing a student with the most suitable feedback for his/her personality, the performed task and environment. The issues of (1) what can be personalized in the feedback and (2) to which characteristics should feedback be personalized are essential in the development of personalized feedback [13].

Response certitude (also called response *confidence* or response *certainty*) specifies the student's certainty in the answer and helps in understanding the learning behavior. The traditional scheme of multiple-choice tests evaluation, where the responses are being treated as absolutely correct or absolutely wrong, ignores the obvious situations when the correct response can be the result of a random or an intuitive guess and luck, and an incorrect answer can be given due to a careless mistake or due to some misconceptions the student may have. Such mistakes are especially crucial in the online assessment, where the evaluation of students' real knowledge and determining students' misconceptions become an even more difficult task for the teacher than in traditional in-class settings. Not allowing for discrimination of these situations may diminish the effects of personalized assessment.

The use of feedback in certitude-based assessment in traditional education has been actively researched for over 30 years [6, 7]. The researchers examined the student's level of confidence in each of the answers and analyzed (1) the differences in performance of students (not) receiving immediate/delayed feedback; (2) how much time a student spent on processing EF; (3) efficiency of feedback in confidence based assessment.

In our earlier pilot experiment and more recently a series of real online assessment studies in [10, 11, 12] we have been able to demonstrate that knowledge of response certitude together with response correctness allows to determine what kind of feedback is more preferable and more effective for the students, and EF may sufficiently improve the performance of students during the online tests. These encouraging results motivated us to develop a recommendation approach for tailoring immediate EF for students' needs in [12]. We presented empirical evidence in [12] that many students are eager to follow the recommendations on necessity or usefulness to read certain EF in the majority of cases, after following the recommendations some students were willing to state explicitly whether particular EF indeed was useful to understand the subject matter better or not (and in most of the cases it was found helpful), and last but not least recommended EF helped to answer related questions better.

Individual LS are one of the important characteristics of the student that characterize the ways in which the student perceives information, acquires knowledge, and communicates with the teacher and with other students. Incorporating LS in WBLSSs has been one of the topical problems of WBLSS design during recent years. There are currently several WBLSSs that support adaptation to LS (AHA!, CS383, IDEAL, MAS-PLANG, INSPIRE). However, according to our knowledge, there is no system or reported research (in the e-learning context) that addressed the issue aimed at providing feedback tailored to the LS of the student except our own recent study [1].

3 Adaptive Selection and Recommendation of Immediate EF

3.1 Authoring Adaptive EF

Feedback adaptation can be based on the traditional user modeling approach in adaptive hypermedia [1]. One key component here is a feedback adaptation unit that has to include a knowledge base containing the adaptation rules that associate user (task, environment) characteristics with certain feedback parameters from the feedback repository. For this particular study we used a simple user model that includes information about student's LS, and certitude and correctness of the current response (which constitute two dimensions of possible cases; high-confidence correct responses (HCCR), high-confidence wrong responses (HCWR), low-confidence correct responses (LCCR), low-confidence wrong responses (LCWR)). Other individual characteristics can be added easily of course, however we tried to focus our study on a particular set of characteristics that allows us to verify our findings from previous experiments as well as to verify the feasibility of the EF adaptation approaches and to make some new observations.

We have studied different aspects of feedback tailoring during a series of experiments (preceding this study) in the form of eight online multiple-choice tests in the Moodle learning system organized as a complimentary yet integral part of three courses (with traditional in-class lectures and instructions) at the Eindhoven University of Technology, the Netherlands during the academic year 2007-2008. Our findings resulted in the implementation of 72 non-contradicting adaptation rules for two types of immediate EF: example-based and theory-based. The base of these rules is compactly summarized in Table 1 below. In the first column, the two dimensions of LS are presented: <[active][balanced][reflective]/[sensing] [balanced][intuitive]>. Cells in the other columns tell what will be directly shown or recommended (number of stars * in the brackets denote the strength of the recommendation) to a student upon the EF request.

3.2 Experiment Design

The online assessment (partial exam) of 73 students of Human-Computer Interaction (HCI) course was organized in March 2008. As in some of the earlier assessments we used feedback adaptation strategies based on student's response correctness and response certitude, and LS.

The online test consisted of 15 multiple-choice questions. The questions were aimed at assessing the knowledge of the concepts and the development of the necessary skills (like understanding of the basic usability rules and problems such as consistency, mapping (between interface and real world), response time problem, etc.). For each answer students had to provide their certitude (which affected the grade) and had a possibility to request and examine EF that could potentially help to answer the related (later) questions better.

Students were not provided with knowledge of (correct) response separately, but they had to infer it from EF instead (if case they were eager to do so). That is the students had to read the explanations of the EF to understand whether their answer

Table 1. The base for adaptation rules

LS	HCCR		LCCR		LCWR		HCWR	
	Show:	Recom- mend:	Show:	Recommend:	Show:	Recommend:	Show:	Recommend:
No L/S	-	-	-	Theory (*) Example (*)	Theory	Example (*)	Theory	Example (***)
Active/ Balanced	-	-	-	Example(**)	Example	Theory(*)	Example	Theory (**)
Reflective/ Balanced	-	Theory (*)	Theory	Example(*)	Theory	Example (**)	Theory	Example(***)
Balanced/ Sensing	-	-	-	Example(**)	Example	-	Example	Theory(**)
Balanced/ Intuitive	-	-	-	Theory(**)	Theory	-	Theory	Example(**)
Active/ Sensing	-	-	-	Example(**)	Example	-	Example	Theory(**)
Active/ Intuitive	-	-	-	Theory (**), Example(*)	Theory	Example (*)	Theory	Example (**)
Reflective/ Sensing	-	Example (*)	-	Example (**) Theory (*)	Example	Theory (**)	Example	Theory (***)
Reflective/ Intuitive	-	Theory (*)	Theory	Example (*)	Theory	-	Theory	Example (***)
Balanced/ Balanced	-	-	-	Theory(*) Example (*)	Theory	Example (*)	Theory	Example (**)

was correct or not. The results of our previous experiments suggested that it is beneficial for the students to embed KR into EF to increase the overall effect of EF on learning process during the assessment.

For every student and for each question in the test we collected all the possible information, including (besides the actual selected answer) correctness, certitude, grade (determined by correctness and certitude), time spent for answering the question, whether feedback was requested or not, and (if it was) which feedback was shown directly, which was recommended with which strength, and finally which one(s) were actually examined (including time spent for examining two each type of feedback in seconds).

Before passing the actual tests the students were asked to complete (not compulsory) Felder-Silverman’s LS quiz (44 questions) [2]; 66 out of 73 students completed this questionnaire.

Adaptation of presentation and recommendation of feedback varied between the questions in the test used for this study. For questions 1, 3, 5, 7, 9, 10, 13, 15 presentation and recommendation of EF was based on student’s LS (active/reflective and sensing/intuitive dimensions), response correctness and response certitude. For the other questions adaptation was performed based only on the response correctness and certitude. For those (few) students who did not complete the (non-mandatory) LS quiz, EF presentation/recommendation was based only on their response correctness and certitude for both groups of questions.

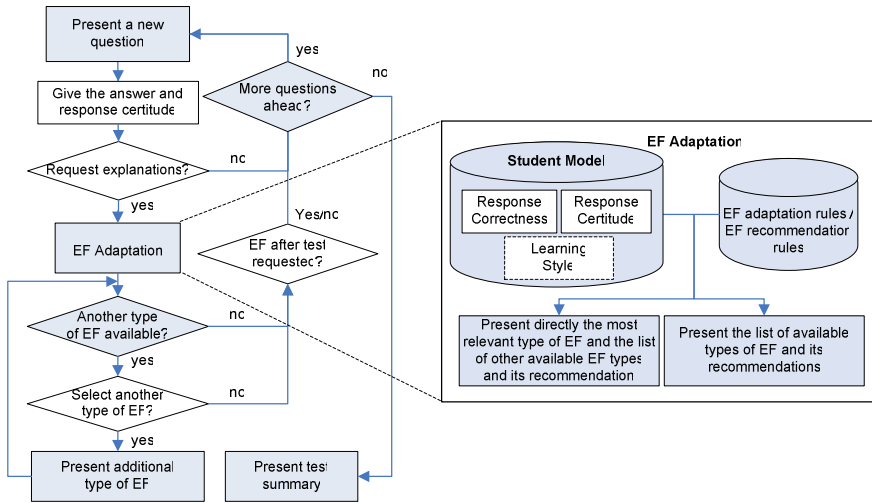


Fig. 1. Assessment process

Further (less important) details regarding the organization of the test, including an illustrative example of the questions and elaborated feedback, are made available in an appendix placed online at <http://www.wis.win.tue.nl/~debra/ectel08/>. Here we only present the flow chart of the assessment process (Fig. 1).

4 Results Obtained

We evaluated the effectiveness of adaptive selection and recommendation comparing the number of requests for the first EF (only in the cases where that EF was not already automatically shown as a result of the adaptation rules and thus did not have to be requested first) and the second EF; the time students spent for studying the adaptively selected or recommended EF (reading vs. scanning the EF); and usefulness of the EF according the students’ feedback rating they provided. The results of earlier experiments already demonstrated that EF sufficiently *improves* the students’ performance during the test. Here we analyze the students’ *perception* of the EF usefulness.

In order to compare two personalization strategies (that is the focus of our analysis here) we analyzed the data from 47 of 73 students for 14 questions (the last question was excluded as an “outlier” in a sense that reading feedback can not help to answer other questions any more from the one hand and on the other hand students should not care about the time limit any longer at this point). We excluded from analysis data also the data of the 7 students who did not complete the LS questionnaire before the test, as for them the personalization/recommendation of EF worked identically for both groups of questions. We also ignored the data of 18 students whose LS was balanced according both dimensions used in personalization (active/reflective, sensing/intuitive), as adaptation rules used in such cases were the same as for personalization based only on response correctness and certitude.

Analysis of the EF requests. Figure 1 illustrates how different EF request-related situation occurred for the questions from Group 1, where adaptive EF selection and recommendation were based on two dimensions of LS (active/reflective, sensing/intuitive) besides the response correctness and certainty, and for the questions from Group 2, where EF adaptation was based only on response certainty and correctness. There were almost equal percentages of initial EF requests in both groups (79% vs. 75%) as well as requests for the explanations in the case no type of EF was directly shown (without the need to request it explicitly): 88% vs. 87,5%. The percentage of requests for additional feedback for Group 2 was higher than for group 1 (27% vs. 16%). This can mean that EF that was shown directly in Group 2 (EF personalization based on response correctness and response certainty) was not always suitable for the students, whereas for the questions from group 1 the type of directly shown feedback was (on average) more suitable for the certain students. Figure 1 also presents the distribution of the responses according to their correctness and certainty (HCCR, LCCR, LCWR, and, HCWR). It helps more clearly to see what the responses of the students were within and between the groups and to analyze how EF adaptation functioned in each case.

For a more detailed comparison of the two EF recommendation/personalization strategies we examine the two most interesting situations: (1) when EF was directly shown to the students (Figure 3 a, b) and (2) when EF was not directly shown, but the user could request one or two available types of EF (Figure 4 a, b).

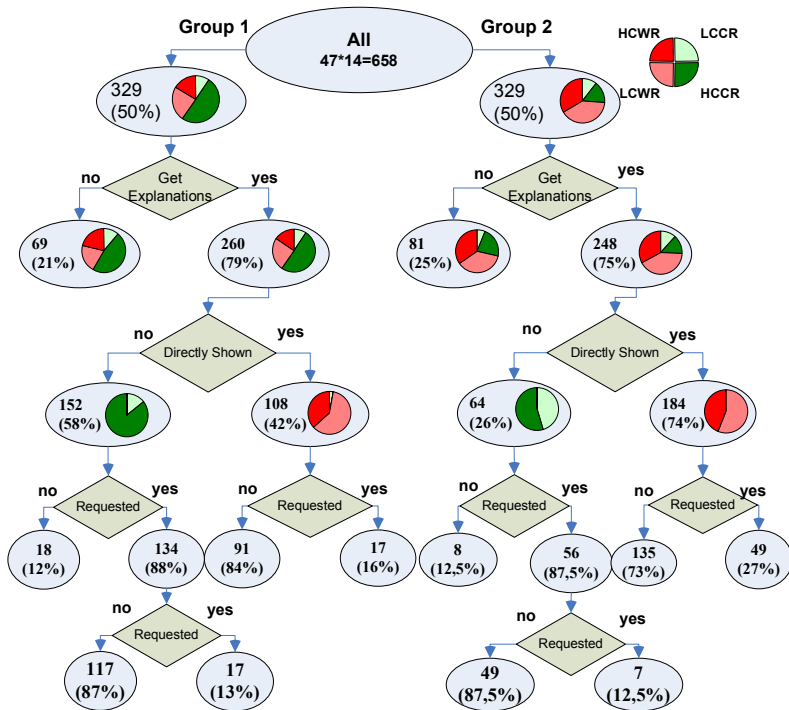


Fig. 2. EF requests related statistics for two groups of questions: Group1 (adaptation rules use LS information) and Group 2 (LS information is not used in adaptation)

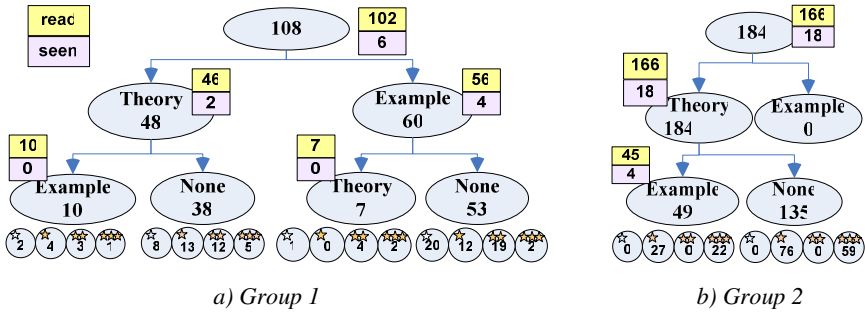


Fig. 3. Student behaviour when one type of EF has been shown directly

In Figure 3 (a,b) we can see that the number of cases where the student was just “scanning” the directly shown EF (marked as “seen” in the figure) was higher (18/166, i.e. 10.8%) for the questions from Group 2 than for the questions from Group 1 (6/102, i.e. 5.9%). This also suggests that the type of directly shown EF was more appropriate for the questions from Group 1 than from Group 2. The percentage of requests for the additional feedback after getting directly shown feedback is also higher for Group 2 with 26.6% (49/(149+35)) vs. 15.7% (17/(53+38+17)) for Group 1. The analysis of the recommendation strength of EF types (that students did or did not request after getting directly shown EF) illustrates that the students followed our recommendations quite well for both groups of the questions. The students requested another available type of EF more often when it was more strongly recommended (with higher number of stars).

In Figure 4 a, b we illustrate the situations when EF was not directly shown, but the students had a possibility to choose it from the two available types of explanations (theory-based and example-based) by either following our recommendations or not. The percentages of requesting theory-based (36% vs. 35%) and example-based (64% vs. 65%) EF were very close (difference is not statistically significant) for Group 1 and Group 2.

In order to measure the quality (or appropriateness) of the recommendation strategies we calculated the corresponding scores as sums of differences between the strength of the recommendation of the requested type of EF and the strength of the recommendation of another available type of EF. The positive coefficient demonstrates that the recommendation strengths of the selected EF were in most of the cases higher than the recommendation strengths of the other available type of EF. For the Group 2 the recommendation of both theory-based and example-based EF were given the same number of stars in the cases where the EF was not shown directly. Thus the calculated scores are illustrative only for the Group 1 in this context. However, for Group 1 the score is positive both for the request of theory-based and example-based types of EF. It can be also seen from the figures that for theory-based EF recommendation with different strengths (blue circles below “Theory 45” in Figure 3a) and recommendations of the example-based EF that was given in those situations (yellow

circle below “Theory 45 in Figure 3a) that the students did request feedback for which the strength of the recommendation was higher. Students requested the second type of feedback after reading or “scanning” the first with the same frequency in both groups (in 12.5% of cases) and followed the recommendation for the another type of feedback available also reasonably well. In general, the score for the next level should be negative, meaning that we can expect that the student, after examining the first selected type of EF, would proceed directly to the next question (if the selected type was suitable). However, in one case (when example-based EF was requested after theory-based) this score was positive which indicates that students often believed that theory-based EF is not adequate or not clear (despite of its recommendation) and hoped that the example-based EF would shed more light on the subject matter.

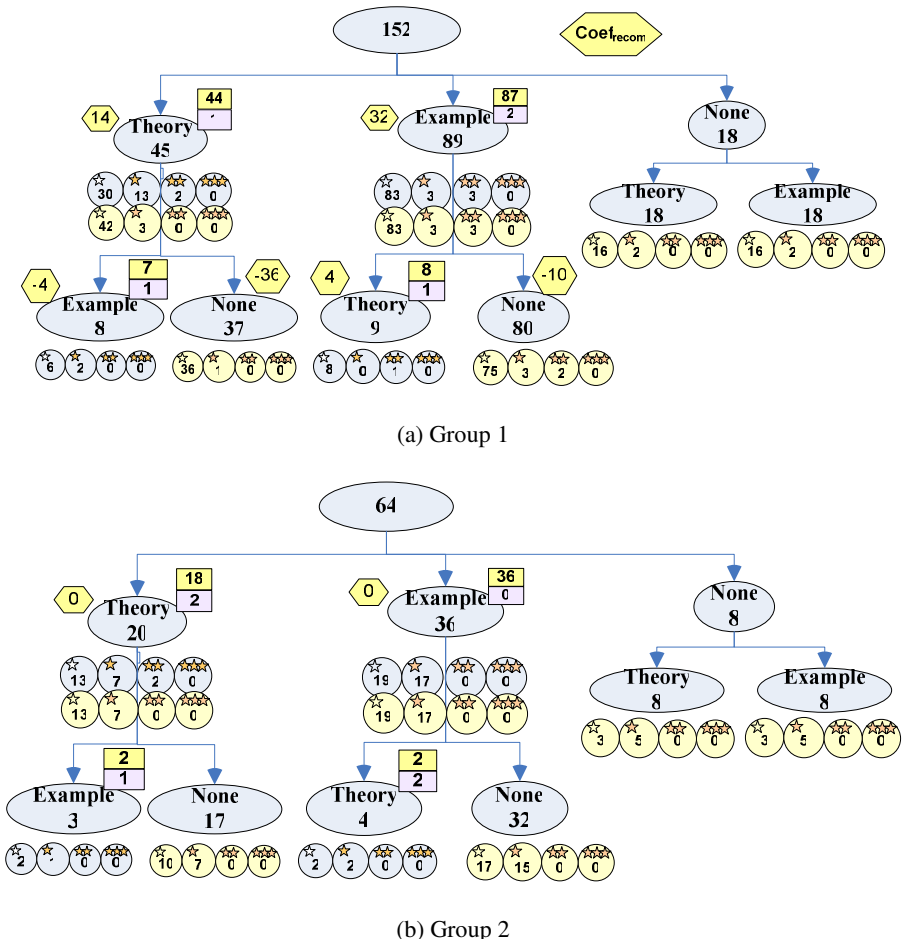


Fig. 4. Student preferences in EF requests when no EF has been shown directly

Usefulness of EF. We analyzed the students’ remarks about the usefulness of EF they were willing to provide during the test. 47 students left in total 82 remarks about EF usefulness; 39 for Group 1 and 43 for Group 2 (providing them was optional; 82 remarks correspond to 14% response rate with respect to the number of actually received EFs by the students). Surprisingly, in both groups in about 18% the students marked EF as not useful, and as useful respectively in 82% (i.e. percentages are almost the same). It is worth noticing that the remarks about not usefulness of EF together with other comments the students provided about the questions and the EF (as a free text typed in the designated places) are taken into account by the teacher for a possible improvement of the test (next year), and also for detecting possible confusion about the questions or answers. The free text comments about the questions and EF were taken into account in the manual re-grading in a few cases.

We also analyzed the recommendation strengths of the EF that the students found useful or not useful. The average recommendation strengths for the EF found to be not useful are higher than for the EF that was found useful. This contradicts an intuition that useful EF is expected to correspond to higher recommendation strengths, but this can be explained by the fact that students tried to provide their evaluation of such EF that was highly recommended but appeared to be not useful (according to the student’s belief). Interestingly also, the ratios of these scores between usefulness/not usefulness for example- and theory-based feedback are very different in Group 1 and Group2 (see Table 2).

Table 2. Average strength of recommendation of EF marked as useful or not useful and percents of students’ remarks about (not) usefulness for 2 groups of questions

		Group 1			Group 2		
		<i>Theory</i>	<i>Example</i>	<i>Total</i>	<i>Theory</i>	<i>Example</i>	<i>Total</i>
Avg. strength of EF recomm.	Useful	0.7	1.6	1.15	3.3	0.6	1.95
	Not useful	2.25	2.7	2.48	4	1.5	2.75
% of students’ remarks	Useful	63.6	89.3	82	86.2	71.4	81.4
	Not useful	36.4	10.7	18	13.8	28.6	18.6

Besides the analysis of how students perceived the usefulness of the different types of EF, we estimated whether EF was helpful in answering related questions students answered. First, we estimated what the relative difference in the performance (grades G) of students is, i.e. the ratio of how many times a “hinted” question $k+c$ was answered better than the question k that contained “hinting” feedback by the students who read that feedback (m students in total) vs. those who did not (n students in total):

$$\sum_{i=1}^m G_{i,k+c} - G_{i,k} / m \quad \text{vs.} \quad \sum_{j=1}^n G_{j,k+c} - G_{j,k} / n$$

Although the relative improvement in Group 1 was more than twice as high (and the difference is statistically significant, $p < 0.05$) as in Group 2 we can not make any

strong conclusions about the advantage of the first adaptation strategy over the second one in this context, because the absolute average improvement of the correctness and grade were rather low (less than 10% for Group 1).

Instead of the direct measurement of grade improvement within the groups as shown above, we also applied several data mining techniques [14], including classification, clustering and association analysis for finding additional evidence of EF usefulness. Mining assessment data appears to be a non-trivial task due to the high inherited redundancy (e.g. grade is identified by correctness and certainty; feedback adaptation/recommendation is defined by the set of rules which use response correctness and certainty and LS) and correlation between the attributes within groups and across the groups (e.g. due to the correlations between the questions). However, it was possible to find some patterns that provide indications of EF usefulness [11].

5 Conclusions and Further Work

Designing and authoring feedback and tailoring it to students is an important problem of online learning assessment. We have studied this problem through a series of experiments in the form of different online tests organized as part of four TU/e courses with traditional in-class lectures and instructions.

In this paper we focused on the immediate EF adaptation by means of adaptive selection and personalized recommendation of the appropriate type of EF for each question answered by the students. Adaptation rules that take into account students' response certitude, response correctness, and LS were designed according to the EF effectiveness and students' preference patterns observed during the preceding studies.

We implemented two adaptation strategies; the first strategy is based on the analysis of response correctness and response certitude only, while the second strategy, besides the analysis of the response, takes student LS into account as well.

Our experimental study demonstrated the feasibility and effectiveness of EF adaptation strategies. The results of the assessment data analysis and as well as feedback received from the students provide enough evidence that our EF adaptation strategies are feasible. In particular, the students (1) followed our recommendations of the type of EF they could select in most of the cases; (2) more often skipped careful examination of EF when it was not directly shown to them as well as EF which they chose by disregarding the recommendations; (3) gave sufficiently more positive than negative responses about the EF that was shown directly or that was recommended to them. According to each of the analyzed dimensions the results obtained either favor (more or less) or at least do not disfavor the second strategy and thus advocate the benefits of taking into account LS for selecting and recommending the most appropriate type of EF during the online assessment.

Our future work on feedback adaptation will be focused on the organization of the similar online assessment studies with more controlled settings for confirming our findings.

Acknowledgments. This research is partly supported by the COMAS Graduate School of the University of Jyväskylä, Finland. We are thankful to the students who participated in the course and provided their valuable comments regarding the organization of the online test and, particularly, the feedback usefulness.

References

1. Brusilovsky, P.: Adaptive hypermedia. *User Modelling and User Adapted Interaction* 11(1/2), 87–110 (2001)
2. Felder, R.M., Silverman, L.K.: Learning and teaching styles in engineering education. *J. of Engineering Education* 78(7), 674–681 (1988)
3. Hatie, J., Timperley, H.: The power of feedback. *J. Review of Educational Research* 87(1), 81–112 (2007)
4. Hummel, H.G.K.: Feedback Model to Support Designers of Blended Learning Courses. *International Review of Open and Distance Learning* 7(3) (2006)
5. Hyland, F.: Providing effective support: investigating feedback to distance language learners. *Open Learning* 16(3), 233–247 (2001)
6. Kulhavy, R.W., Stock, W.A.: Feedback in written instruction: The place of response certitude. *Educational Psychology Review* 1(4), 279–308 (1989)
7. Mory, E.H.: Feedback research revisited. In: Jonassen, D. (ed.) *Handbook of research on educational communications and technology*, pp. 745–783. Lawrence Erlbaum, Mahwah (2004)
8. Pechenizkiy, M., Calders, T., Vasilyeva, E., De Bra, P.: Mining the Student Assessment Data: Lessons Drawn from a Small Scale Case Study. In: *Proc. of 1st Int. Conf. on Educational Data Mining EDM* (to appear, 2008)
9. Shute, V.J.: Focus on formative feedback, Research Report (Retrieved January 15, 2008) (2007), <http://www.ets.org/Media/Research/pdf/RR-07-11.pdf>
10. Vasilyeva, E., Pechenizkiy, M., De Bra, P.: Adaptation of Elaborated Feedback in e-Learning. In: *AH 2008. LNCS*, vol. 5149, pp. 235–244. Springer, Heidelberg (2008)
11. Vasilyeva, E., De Bra, P., Pechenizkiy, M., Puuronen, S.: Tailoring feedback in online assessment: influence of learning styles on the feedback preferences and elaborated feedback effectiveness. In: *Proc. of 8th IEEE Int. Conf. on Advanced Learning Technologies ICALT 2008*. IEEE CS Press, Los Alamitos (2008)
12. Vasilyeva, E., Pechenizkiy, M., De Bra, P.: Tailoring of feedback in web-based learning: the role of response certitude in the assessment. In: Woolf, B.P., Aïmeur, E., Nkambou, R., Lajoie, S. (eds.) *ITS 2008. LNCS*, vol. 5091, pp. 771–773. Springer, Heidelberg (2008)
13. Vasilyeva, E., Puuronen, S., Pechenizkiy, M., Räsänen, P.: Feedback adaptation in web-based learning systems. *Special Issue of Int. J. of Continuing Engineering Education and Life-Long Learning* 17(4-5), 337–357 (2007)
14. Witten, I., Frank, E.: *Data Mining: Practical Machine Learning Tools and Techniques*, 2nd edn. Morgan Kaufmann, San Francisco (2005)