# Domain Specific Interactive Data Mining

Roland Hübscher[1], Sadhana Puntambekar[2], and Aiquin H. Nye[3]

[1] Department of Information Design and Corporate Communication, Bentley College, 175 Forest Street, Waltham, MA 02452-4705, U.S.A.
rhubscher@bentley.edu
[2] Department of Educational Psychology, University of Wisconsin, 1025 W Johnson St., Rm 880, Madison, WI 53706, U.S.A.
puntambekar@education.wisc.edu
[3] Raytheon Company, T2SJ01, 50 Apple Hill, Tewksbury, MA 01876, U.S.A.
quin@raytheon.com

**Abstract.** Finding patterns in data collected from interactions with an educational system is only useful if the patterns can be meaningfully interpreted in the context of the student-system interaction. To further increase the chance of finding such meaningful patterns, we extend the mining process with domain and problem specific representations and the pattern detection expertise of qualified users. The user, that is, the researcher looking for patterns, is not just evaluating the result of an automatic data mining process, but is actively involved in the design of new representation and the search for patterns. This approach is used in addition to more traditional methods and has resulted in a deeper understanding of our data.

## 1 Introduction

In the preface to the Educational Data Mining Workshop at ITS 2006, educational data mining is defined as "the process of converting raw data from educational systems to useful information that can be used to inform design decisions and answer research questions" [1]. Information is only useful if it can be meaningfully interpreted in the appropriate context, for instance, in the context of the student-system interaction. Many data and information representations and many mining algorithms exist from which the user,[4] the researcher interested in understanding the data, can choose. It is not uncommon, that the process of developing representations and mining algorithms is separate from mining actual data, done by different groups of researchers taking advantage of their special areas of expertise. However, since knowledge about the problem domain is important to select the appropriate representations and methods, this can also be a disadvantage, especially if the appropriate methods and representations are not readily available. In that case, the user of the mining tool is forced to use whatever is available.

Discovering useful characteristics of data is not a simple method where data is fed into some black box and the interesting characteristics are computed and returned to the user. Mannila, for instance, suggests a process consisting of the following steps: "1. understanding the domain, 2. preparing the data set, 3. discovering patterns (data mining), 4. postprocessing of discovered patterns, and 5. putting the results into use" [2]. Based on our search for informative patterns in our data, we suggest a similar process. However, we emphasize its iterative nature based on a design process and we will describe and illustrate the specific steps with a concrete example from our own mining efforts. Furthermore, although other researchers may implicitly use a similar approach [3–5], it is important that the process is made explicit so that it can be discussed, shared, improved and followed.

Our educational hypermedia system CoMPASS uses dynamic concept maps to support navigation. We are interested in using the logged navigation data to understand how the student-computer interaction can be related to the student's learning strategies and understanding of the subject matter presented by the system. We intend to use the found relationships between student behavior and student learning to provide adaptive prompts to scaffold the learner as well as to provide teachers with realtime feedback about the students performance [6].

---

[4] We use the following terminology in this paper. The *user* is the person interested in finding patterns. The *learner* or *student* is the person using the educational system.

We only can accept data mining results that can be interpreted meaningfully in the context of the learner using CoMPASS with its specific interface and structure. Sometimes, "interesting" relationships can be found, yet mapping them meaningfully back into the domain where a learner is interacting with a specific system proves very difficult and sometimes even impossible. Thus, we have adopted a method that allows the researcher looking for meaningful patterns in the log data to be part of the mining process. The mining process is an interaction between computer and researcher, both helping each other to find the relationships between log data and student behavior.

This interactive process does not seem to be the norm. Some definitions suggest that the mining process is automatic, for instance, Wikipedia defines data mining as "the process of automatically searching large volumes of data for patterns using tools such as classification, association rule mining, clustering, etc." [7]. This suggests that the researcher interested in the potential patterns in the raw data is not really involved in the data mining process, but only in the interpretation of whatever the data mining algorithm produces. We propose to use a less narrow view.

The goal of this paper is to discuss the advantages and disadvantages of

- involving the user at various stages in the pattern discovery process,
- the use of domain and problem specific algorithms and representations, and
- the use of an iterative design and discovery process.

In this paper, we describe the interactive mining process we have been using to make sense of the raw data collected from the use of CoMPASS. We illustrate this general method with the domain-specific algorithms and representations used to mine our data for meaningful patterns. We will focus on the method, but will not address in detail the specific algorithms or the insights we have gained in the pedagogical domain. For some of these other results, see [8].

## 2    Interactive Data Mining

Interactive data mining allows the user and the the data mining algorithms to interact with each other. Often, the data is visualized helping the user to understand the patterns better and also allowing the user, and not just the mining tool's discovery engine, to discover some of the patterns [9, 10]. Efforts to build integrated environments like VDM [11] supporting many data mining and visualization techniques are of great value and we hope, that at some time in the future, we will be able to extend such a system in the way described here. While a tool like VDM gets its power through the many different mining and visualization techniques it provides so that they can be applied almost effortlessly in many domains with different data, we are interested in enabling tools to add specialized algorithms and visualizations relatively easily with some end-user programming tool. This is a long term goal. For now, we simulate this with a set of programs written as needed in the flexible programming language Python.

Our process is based on our work on log data collected from students interacting with CoMPASS. Before we discuss the specifics of CoMPASS and how we analyzed that log data, we present the process in a more general way and then address each step separately.

1. Collect raw data from learner-system interaction
2. Analyze system and its users, use and context
3. Represent raw data in a meaningful way using domain-specific methods
4. Find clusters of similar data points
5. Visualize (members of) clusters
6. Interpret visualizations in the context of the learner-system interaction
7. If results are not good enough (and we have more time), go back to step 2

This process has some similarities with an iterative design process [12] where the understanding of the problem co-evolves with the solution. In other words, as the user is mining the data, the user learns more about the data and will be able to find more appropriate representations and methods. This does indeed require the mining tool to be extensible with some end-user programming language. It also requires the user to be aware that data mining is not a one-shot approach. Initial results need to be used to improve the mining approach to find more interesting results.

We focus on clustering and ignore some other useful methods. Although we do not want to exclude all other methods, clustering allows us to include the user as part of the mining process in a relatively straightforward manner.

Let's start with the first step of the process. Of course, first the raw data has to be collected. It is important that the data is always analyzed with the context in mind in which the data was generated. Thus, it would be a big mistake to collect the data and then hand it off to a data miner who is completely unaware of the learners' characteristics, the educational system and other factors influencing the learner-system interaction and expect that the miner would return anything terribly meaningful. In other words, the patterns are not just in the data. After all, as soon as we talk about patterns there is a bias[5] involved. Since we cannot avoid some bias completely, it should at least be a result of our understanding of the learner-system interaction including the system interface, the learner characteristics and the pedagogical methods used.

Most of the time, we probably do not want to cluster the raw data, but a more meaningful description of it. How the data should be represented depends on the specific circumstances, of what answers the user wants to answer and the data itself. For instance, as we shall show in the next section, we were not so much interested in finding similar behavior, but in finding similar understanding. Thus, we represented the data so that it would capture more the students' understanding than just their behavior.

Clustering the data requires that we develop some kind of distance or similarity measure further biasing the whole discovery process. We again propose to use domain specific metrics that are consistent with the represented data and the questions the user wants to answer with the analysis.

So far, the user has been involved in the process by selecting representations and similarity metrics or possibly developing them anew based on the understanding of the data and the patterns already found in earlier iterations. Once the clusters have been found, the user has to decide how to analyze their members to find common characteristics or patterns. Since we propose to put the burden for finding patterns, at least to some degree, on the user, visualizations may be useful here. And again, domain-specific visualizations should be considered, although standard ones should be used if they are adequate for the current situation. Just creating domain-specific representations for their own sake is a bad idea since developing them is very time consuming.

When the user studies these clusters for interesting patterns, it is important that the interpretation of these patterns must be done within the context in which the raw data was collected. Thus, the circle closes and the user should go back and consider modifying or completely changing some of the representations used based on what was learned during the previous iteration. As the user iterates through the process, the understanding of the data, patterns and their meaning evolves. Finding the answers is not a one-shot approach.

Our proposed method is also somewhat analogical to how expert systems have evolved over the last thirty years. The early expert systems used to ask the user for various inputs, do some reasoning and return the result, or a list of results with some associated confidence factors. Although that mode of operation can be useful under certain circumstances, intelligent systems are viewed now more and more as intelligent assistants helping the user solve the problem collaboratively [13].

## 3   Mining CoMPASS' Navigation Logs

We now illustrate the ideas introduced in the previous section with the data analysis of the navigation data collected with CoMPASS. CoMPASS is an educational hypermedia system with navigation support in the form of dynamic concept maps [8]. CoMPASS helps students understand the relationships between science concepts and principles. It uses two representations, concept maps and text, to support navigation and learning. Each page in CoMPASS represents a conceptual unit such as force or acceleration. A conceptual map of the science concept and other related concepts takes up the left half of the CoMPASS screen, and a textual description takes up the right half of the screen (see Figure 1). The maps are dynamically constructed and displayed with the fisheye technique every time the student selects a concept. The selected (focal) concept is at the center of the map, with the most related concepts at the first level of magnification and those

---

[5] We use bias in the non-technical sense throughout this paper.

less closely related at the outer level of the map. The maps in CoMPASS mirror the structure of the domain to aid deep learning and are designed to help students make connections, giving students alternative paths to pursue for any particular activity, so that they can see how different phenomena are related to each other.



**Fig. 1.** CoMPASS with navigation support on the left and a description of the concept force in the context of an inclined plane.

We are interested in understanding the navigation paths of the students in CoMPASS for several reasons. The nonlinear nature of hypertext can be used to organize information in multiple ways, reflecting the structure of the described domain. As a result, navigation through hypertext requires the learner to make frequent decisions and evaluations. Providing the proper navigation support is therefore important and understanding how navigation and learning relates to each other is therefore important. Furthermore, we intend to provide adaptive support to the students in form of dynamic prompts triggering metacognitive activities. Such prompts have to be sensitive to the learning context including the students' understanding and potential problems. We hope that we can associate certain navigation patterns with students' understanding to provide the adequate prompts in real time. Similar to [14], we also want to detect in real time students who may have some learning problems so that the teacher, a highly valuable but sparse classroom resource, can focus his or her attention on those students who need it the most.

Before we discuss the steps introduced in the previous section, it is important that we make the questions we are interested in with respect to the data logged in CoMPASS explicit. This allows us to develop the domain and problem specific representations with the research questions and the learner-CoMPASS interaction in mind. One of the goals of CoMPASS is, together with other class room interventions, to scaffold students to gain a deep understanding of the domain specific concepts and their relationships. In other words, we are interested in the students' structural (or relational) knowledge [15]. In the case of the content displayed in Figure 1, the topics are simple machines (e.g., inclined plane, lever, screw) and the concepts are from the domain of mechanics and include energy, force, efficiency and gravity as the concept map shows.

In CoMPASS, navigation data is collected in the form of a sequence of navigation events. Each event consists of the time of the mouse click, the name of the student who clicked on it

and the destination page. Since each page contains the description of exactly one concept, every destination page is equivalent to a destination concept. This is not a very rich data source and we were initially worried that we might not find interesting patterns. In addition, the individual interactions are relatively short, that is, the students rarely click on more than twenty links in one session whose duration is normally between 60 and 90 minutes. For each user, the raw data is then collected in an $n \times n$ navigation matrix $N$ such that $N_{ij}$ is the number of transitions from concept $i$ to concept $j$. A transition from $i$ to $j$ simply means that the user, being on the page for concept $i$, has clicked on a link to the page describing concept $j$.

The next step requires to represent the raw data $N$ to increase the chance of finding patterns that address the questions we are interested in. In other word, the new representation needs to have characteristics we consider to be relevant in similar students. Since we are interested in the structural knowledge of a student, we wanted a representation that would emphasize the structure hidden within the navigation data. For this purpose, we applied the Pathfinder Network Scaling procedure computing an approximate representation of the conceptual model of the user [16]. The Pathfinder algorithm was developed to find relevant relations in data that describes the proximity between concepts. Naturally, all concepts are somehow related to all others, however, only the relevant relations should be retained. The Pathfinder algorithm has been successfully used for this task in various domains [17]. We modified the algorithms so that it works for navigation networks where two concepts are closer if there are more direct transitions between them. The resulting Pathfinder network is again an $n \times n$ matrix and can be interpreted as a concept map representing the structural knowledge of an individual learner (see Figure 2 for an example).
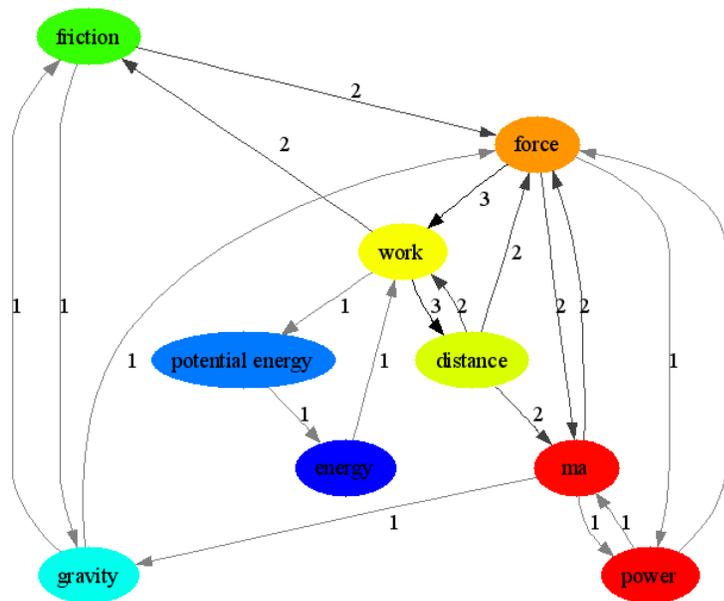


**Fig. 2.** The output of the Pathfinder algorithm which can be interpreted as the concept map describing a student's structural knowledge.

Before the user can look for patterns in the data, the data points need to be clustered [18]. In our case, these data points are the learner models, that is, the Pathfinder networks. We originally applied the k-Means clustering algorithm [19] because of its simplicity and adequate results. Clustering requires some function that measures the similarity (or distance) between two data points. Again, we chose one that was consistent with our interest in the structural characteristics of the learner models. After some testing, we settled on a simple measure suggested by the inventor of the Pathfinder methods [16] which measures the structural similarity of graphs, that is, the Pathfinder networks representing the students' understanding.

Given are two Pathfinder networks $P$ and $Q$ and we want to compute their structural similarity $sim(P, Q)$. We can assume that they have the same size $n \times n$ and that their node labels are ordered

the same in both graphs. If that's not the case, we simply extend both graphs to include all labels and order them lexicographically. However, we do not include any nodes that neither network connects to.

Let $P_{ij}$ and $Q_{ij}$ be the vertex from $i$ to $j$ in $P$ and $Q$, respectively. Since the vertices are ordered, the indices refer to the vertices with the same labels in both networks. Then, the similarity is computed by averaging over the structural similarity of all vertices. The similarity of vertex $i$ in $P$ and vertex $i$ in $Q$ is the the intersection of vertex $i$'s respective outgoing edges divided by the union of the same edges. Since the edges are weighted by the number of transitions, union and intersection are computed as the maximum and minimum, respectively, of the edges' weights. Since the

$$\text{sim}(P,Q) = \frac{1}{n} \sum_{i=1}^{n} \frac{\sum \min_{j=1}^{n}(P_{ij}, Q_{ij})}{\sum \max_{j=1}^{n}(P_{ij}, Q_{ij})}$$

Although, the results we obtained with k-Means were satisfactory, though somewhat unstable—like many other greedy algorithms, k-Means does not always find the optimal solution—we have also used hierarchical clustering which provides more fine-grained information [20]. K-Means clustering creates a partition of the learner models which then are visualized as discussed below. However, when using hierarchical clustering, it is possible to look at many more meaningful sub-groups depending on where the cutoff is made as Figure 3 shows. In this figure, the names on the left are the names of the learners. It shows that students *green3* and *red3* are quite similar and so are *purple3* and *yellow3*. So these two clusters can be visualized to see what they have in common, but also the visualizations for the cluster consisting of all four students is generated and so on. As the dendrogram in Figure 3 shows, five meaningful clusters and subclusters are generated and can be visualized.
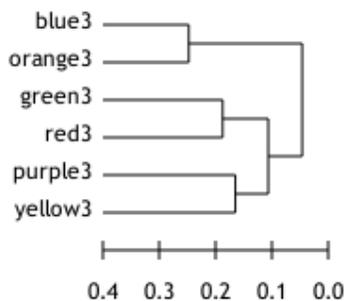


**Fig. 3.** The hierarchical clusterer computes a dendrogram as output.

In the k-Means and the hierarchical clustering algorithm we used the centroid distance function where the distance between two clusters is measured by the distance between the centroids of the two clusters. The centroid of a cluster is the average of all the data points in that clusters, in our case, the average of the Pathfinder matrices.

The next step is, as already mentioned, visualizing the clustered networks. We visualize all clusters in a hierarchical clustering for further study. However, once the similarity becomes small, finding interesting patterns tends to becomes less probable, because the accumulation of several not so similar learner models results in a "washout" effect: in average, each concept is a bit related to each other and nothing characteristic stands out. This, for instance, tends to be true for the trivial cluster including all of the students.

These clusters serve as the starting point for the user to find interesting relationships. Instead of visualizing these clusters in some standard form—we do that, too—we put much effort into finding visualizations that are meaningful with respect to how the students use CoMPASS and how CoMPASS is structured. One obvious representation is the accumulated models, that is, we average the network outputs by Pathfinder for all students in the cluster which results in

a network similar to the one show in Figure 2, however, as mentioned the washout effect is a problem.

Before we turned to the type of visualizations described below, we studied the centroids of the clusters like the one in Figure 4. We did indeed find interesting patterns and were able to relate them to the students' learning [8]. Some students were rather focused and explored more or less other topics, others showed a random "pattern" and the ones in Figure 4 a highly linear behavior influenced by the interface. Random and linear behaviors correlated with relatively low learning performance. Although this analysis was quite successful, we are interested in finding additional less obvious patterns with visualizations that hopefully make these patterns easier to recognize.



**Fig. 4.** We also analyzed centroids of the clusters.

Examples of visualizations that are much more domain specific are shown in Figures 5 and 6. Instead of providing an aggregate view for a cluster, each cluster member is displayed separately in form of a ring graph (see Figure 5). The ring is based on some important characteristics of CoMPASS and its use as explained below.



**Fig. 5.** A ring graph describing what descriptions students visited during a session. The outer ring refers to concepts in the context of a topic, the inner ring to context-free definitions.

The domain-specific visualizations are being used by the researchers familiar with CoMPASS and its use. Thus, to understand ring graphs as used here, some details of CoMPASS need to be further explained. CoMPASS provides various types of concept descriptions for middle school students. The types refer to the context in which the concepts are described. For instance, the concept of force can be described in the context of falling objects or in the context of an inclined plane as in Figure 1. In CoMPASS, a concept description without context is called the concept's definition. Since we consider the distinction between descriptions within and without a context pedagogically meaningful in CoMPASS, the two concentric rings in Figure 5 capture this characteristic of CoMPASS. The ring represents a session of using CoMPASS starting at the top and moving clockwise around the ring. The inner ring represents visits by the student to definitions, the outer ring visits to descriptions within some context. Different colors are used to code what concepts are described. We found that it was relatively easy to pick up meaningful patterns by people familiar with CoMPASS and the student-system interaction in which the data is collected.

A new representation we have been working on is the panel graph in Figure 6 where different students are represented with different colors. There are six sections from left to right. The left-most section refers to definitions (context free), the next one to concept descriptions in the context of inclined plane, then in the context of lever, and so on. The navigation events are ordered starting at the top of the graph and going down. Again, this representation captures important relationships of CoMPASS and its use and may support finding interesting behavioral patterns.
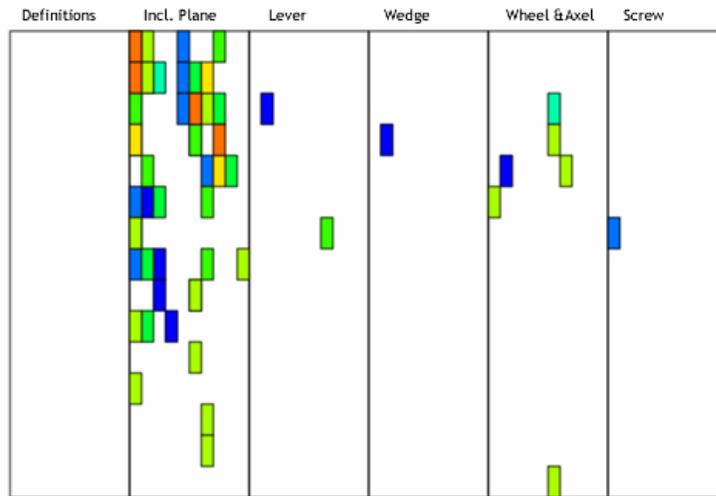


**Fig. 6.** A panel graph comparing the navigation behavior of various student groups. Each group is represented by a specific color.

What is important here is not the ring or panel graph per se, but that it was designed iteratively with the student-system interaction and the research questions in mind. Domain and problem specificity can be quite powerful, though developing these representations is not trivial and takes time. However, having representations that are relatively easy to interpret with respect to the actual research questions makes the representations very useful. Patterns mean immediately something whereas in other situations, patterns are found and then it is sometimes difficult to figure what they actually mean.

## 4 Conclusions

We are interested in finding meaningful patterns in the data collected from the interaction between students and the educational hypermedia system CoMPASS. For instance, we have studied the navigation data also with methods from social network analysis [21] and have found some interesting patterns, however, it has been quite difficult to make sense of these relationships at a

pedagogical level. Just pointing out some interesting commonalities that do not have a meaningful interpretation are simply not useful in general.

Therefore, we have proposed an approach that takes advantage of domain and problem specific knowledge and human experts as pattern finders. We do not imply that all data mining should follow the proposed method, but see it more of a way of using and possibly extending existing tools. Our implementation is at the moment still relatively ad hoc where new domain-specific representations and algorithms have to implemented "by hand" in Python. This is quite costly and it is not obvious that an integrated environment could much more easily be extended with new representations.

We have addressed the reasons for using domain specific representations and visualizations and its advantages. However, this approach also has potential disadvantages. As soon as one makes assumptions about what characteristics are interesting and which ones are not, a bias is introduced which may prevent certain patterns from being found. For instance, our focus on the structural knowledge of the students is justified given our research questions, however, it also may keep certain interesting and meaningful relations hidden. After all, one can only see what one displays and as soon as one emphasizes one property, another is being deemphasized [22]. Therefore, the proposed method should be used in addition to more general approaches, not as their replacement.

## Acknowledgments

## References

1. Heiner, C., Baker, R., Yacef, K.: Preface. In: Workshop on Educational Data Mining at the 8th International Conference on Intelligent Tutoring Systems (ITS 2006), Jhongli, Taiwan. (2006)
2. Mannila, H.: Methods and problems in data mining. In Afrati, F., Kolaitis, P., eds.: International Conference on Database Theory, Delphi, Greece, Springer Verlag (1997) 41–55
3. Kay, J., Maisonneuve, N., Yacef, K., Zaïane, O.: Mining patterns of events in students' teamwork data. In: Proceedings of the Workshop on Educational Data Mining at the 8th International Conference on Intelligent Tutoring Systems (ITS 2006), Jhongli, Taiwan. (2006)
4. Merceron, A., Yacef, K.: Educational data mining: a case study. In: Proceedings of the 12th Conference on Artificial Intelligence in Education, Amsterdam, The Netherlands. (2005)
5. Mazza, R., Dimitrova, V.: Visualising student tracking data to support instructors in web-based distance education. In: Proceedings of the Thirteenth International World Wide Web Conference (WWW2004), New York. (2004)
6. Puntambekar, S.: Analayzing navigation data to design adaptive navigation support in hypertext. In Hoppe, U., Verdejo, F., Kay, J., eds.: Artificial Intelligence in Education: Shaping the future of learning through intelligent technologies, IOS Press (2003) 209–216
7. Wikipedia: Data mining. Retrieved January 27, 2007, from: http://en.wikipedia.org/wiki/Data_mining (2007)
8. Puntambekar, S., Stylianou, A., Hübscher, R.: Improving navigation and learning in hypertext environments with navigable concept maps. Human-Computer Interaction **18**(4) (2003) 395–428
9. Aggarwal, C.C.: Towards effective and interpretable data mining by visual. ACM SIGKDD Explorations Newsletter **3** (2002) 11–22
10. Spenke, M., Beilken, C.: Visual, interactive data mining with infozoom – the financial data set. In: "Discovery Challenge" at the 3rd European Conference on Principles and Practice of Knowledge Discovery in Databases, PKDD 99, Prague, Czech Republic. (1999)
11. Schulz, H.J., Nocke, T., Schumann, H.: A framework for visual data mining of structures. In: Twenty-Ninth Australasian Computer Science Conference (ACSC2006), Hobart, Tasmania. (2006)
12. Nielsen, J.: Iterative user-interface design. Computer **26**(11) (1993) 32–41
13. Hoschka, P.: Computers As Assistants: A New Generation of Support Systems. Lawrence Erlbaum Associates (1996)
14. Ma, Y., Liu, B., Wong, C.K., Yu, P.S., Lee, S.M.: Targeting the right students using data mining. In: Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco. (2001)

15. Jonassen, D.H., Beissner, K., Yacci, M.: Structural Knowledge: Techniques for Representing, Conveying, and Acquiring Structural Knowledge. Lawrence Erlbaum Associates, Hillsdale, NJ (1993)
16. Schvaneveldt, R.W., ed.: Pathfinder Associative Networks: Studies in Knowledge Organization. Ablex, Norwood (1990)
17. Chen, C.: Visualizing semantic spaces and author co-citation networks in digital libraries. Information Processing & Management **35**(3) (1999) 401–420
18. Merceron, A., Yacef, K.: TADA-Ed for educational data mining. Interactive Multimedia Electronic Journal of Computer-Enhanced Learning (IME$_j$) **7**(1) (2005)
19. Hansen, P., Mladenovic, N.: J-Means: A new local search heuristic for minimum sum-of-squares clustering. Pattern Recognition **34**(2) (2001) 405–413
20. Jain, A.K., Murty, M.N., Flynn, P.J.: Data clustering: a review. ACM Computing Surveys **31**(3) (1999) 264–323
21. Wasserman, S., Faust, K.: Social network analysis. Cambridge University Press, New York, NY (1994)
22. Narayanan, N.H., Hübscher, R.: Visual language theory: Towards a human-computer interaction perspective. In Meyer, B., Marriott, K., eds.: Visual Language Theory. Springer Verlag (1998) 85–127