# Information Visualization, Visual Data Mining and Machine Learning

Daniel A. Keim, Fabrice Rossi, Thomas Seidl, Michel Verleysen, and Stefan Wrobel

Information visualization and visual data mining leverage the human visual system to provide insight and understanding of unorganized data. Visualizing data in a way that is appropriate for the user's needs proves essential in a number of situations: getting insights about data before a further more quantitative analysis, presenting data to a user through well-chosen table, graph or other structured representations, relying on the cognitive skills of humans to show them extended information in a compact way, etc.

Machine learning enables computers to automatically discover complex patterns in data and, when examples of such patterns are available, to learn automatically from the examples how to recognize occurrences of those patterns in new data. Machine learning has proven itself quite successful in day to day tasks such as SPAM filtering and optical character recognition.

Both research fields share a focus on data and information, and it might seem at first that the main difference between the two fields is the predominance of visual representations of the data in information visualization compared to its relatively low presence in machine learning. However, it should be noted that visual representations are used in a quite systematic way in machine learning, for instance to summarize predictive performances, i.e., whether a given system is performing well in detecting some pattern. This can be traced back to a long tradition of statistical graphics for instance. Dimensionality reduction is also a major topic in machine learning: one aims here at describing as accurately as possible some data with a small number of variables rather than with their original possibly numerous variables. Principal component analysis is the simplest and most well known example of such a method. In the extreme case where one uses only two or three variables, dimensionality reduction is a form of information visualization as the new variables can be used to directly display the original data.

The main difference between both fields is the role of the user in the data exploration and modeling. The ultimate goal of machine learning is somehow to get rid of the user: everything should be completely automated and done by a computer. While the user could still play a role by, e.g., choosing the data description or the type of algorithm to use, his/her influence should be limited to a strict minimum. In information visualization, a quite opposite point of view is put forward as visual representations are designed to be leveraged by a human to extract knowledge from the data. Patterns are discovered by the user, models are adjusted to the data under user steering, etc.

This major difference in philosophy probably explains why machine learning

and information visualization communities have remained relatively disconnected. Both research fields are mature and well structured around major conferences and journals. There is also a strong tradition of Dagstuhl seminars about both topics. Yet, despite some well known success, collaboration has been scarce among researchers coming from the two fields. Some success stories are the use of state-of-the-art results from one field in the other. For instance, Kohonen's Self Organizing Map, a well known dimensionality reduction technique, has been successful partly because of its visualization capabilities which were inspired by information visualization results. In the opposite direction, information visualization techniques often use classical methods from machine learning, for instance, clustering or multidimensional scaling.

The seminar was organized in this context with the specific goal of bringing together researchers from both communities in order to tighten the loose links between them. To limit the risk of misunderstandings induced by the different backgrounds of researchers from the two communities, the seminar started with introductory talks about both domains. It was then mainly organized as a series of thematic talks with a significant portion of the time dedicated to questions and discussions. After the first two days of meeting, understanding between both communities reached a sufficient level to organize, in addition to the plenary talks, working group focusing on specific issues.

Several research topics emerged from the initial discussions and lead to the creation of the working groups. The subject that raised probably the largest number of questions and discussions is Evaluation. It is not very surprising as differences between the communities about evaluation (or quality assessment) might be considered as the concrete technical manifestation of cultural and philosophical differences between them. Indeed, in machine learning, automatic methods are mostly designed according to the following general principle: Given a quality measure for a possible solution of the problem under study, one devises an algorithm that searches the solution space efficiently for the optimal solution with respect to this measure. For instance, in SPAM filtering a possible quality measure is the classification accuracy of the filter: it has to sort unsolicited bulk messages correctly into the SPAM class and all other emails in the HAM class. In a simple setting, the best filter could be considered as the one with the smallest number of errors. However, counting only the number of errors is usually too naive, and better quality measures have to be used, such as the area under the ROC curve: the Receiver Operating Characteristic curve shows the dependency between the true positive rate (the percentage of unsolicited bulk messages classified as SPAM) and the false positive rate (the percentage of correct emails classified as SPAM).

In information visualization, evaluation cannot rely only on mathematical quality measures as the user is always part of the story. A successful visualization is a solution, with which the user is able to perform better, in a general sense, compared to existing solutions. As in machine learning, a method is therefore evaluated according to some goal and with some quality metric, but the evaluation process and the quality metrics have to take the user into account. For instance, one display can be used to help the user assess the correlation between variables. Then, a quality metric might be the time needed to find a pair of highly correlated variables, or the time needed to decide that there is no such pair. Another metric might be the percentage of accurate decisions about the correlation of some pairs of variables. In general, a visualization system can be evaluated with

respect to numerous tasks and according to various metrics. This should be done in a controlled environment and with different users, to limit the influence of interpersonal variations.

Among the discussions between members of the two communities about evaluation, questions were raised about the so-called unsupervised problems in machine learning. These problems, such as clustering or dimensionality reduction, are ill-posed in a machine learning sense: there is no unquestionable quality metric associated to e.g. clustering but rather a large number of such metrics. Some of those metrics lead to very difficult optimization problems (from a computational point of view) that are addressed via approximate heuristic solutions. In the end, machine learning has produced dozens of clustering methods and dimensionality reduction methods, and evaluations with respect to user needs remain an open problem. An important outcome of the seminar was to reposition this problem in the global picture of collaboration between information visualization and machine learning. For instance, if many quality measures are possible, one way to compare them would be to measure their link to user performances in different tasks. If several methods seem to perform equally well in a machine learning sense, then the user feedback could help to indentify the "best" method. It was also noted that many methods that are studied in machine learning and linked to information visualization, in particular dimensionality reduction and embedding techniques, would benefit from more interaction between the communities. At minimum, state-of-the-art methods from machine learning should be known by information visualization researchers and state-of-the-art visualization techniques should be deployed by machine learning researchers.

Another topic discussed thoroughly at the seminar was the visualization of specific types of objects. Relational data were discussed, for instance, as a general model for heterogeneous complex data as stored in a relational database. Graph visualization techniques provide a possible starting point, but it is clear that for large databases, summarization is needed, which brought back the discussion of the ill defined clustering problem mentioned above. Among complex objects, models obtained by a machine learning algorithms were also considered, in particular as good candidates for interactive visualizations. Decision trees give a good example of such objects: Given a proper visualization of the current tree, of some possible simplified or more complex versions and of the effect of the tree(s) on some dataset, an expert user can adapt the tree to his/her specific goals that are not directly expressible in a quality criterion. The extreme case of visualizing the dynamic evolution of a self learning process was discussed as a prototype of complex objects representation: The system is evolving through time, it learns decision rules, and it evolves using complex (and evolving) decision tables.

Finally, it became clear that a large effort is still needed at the algorithmic and software levels. First, fast machine learning techniques are needed that can be embedded in interactive visualization systems. Second, there is the need for a standard software environment that can be used in both communities. The unavailability of such a system hurts research to some extent as some active system environments in one field do not include even basic facilities from the other. One typical example is the R statistical environment with which a large part of machine learning research is conducted and whose interactive visualization capabilities are limited, in particular in comparison to the state-of-the-art static visualization possibilities. One possible solution foreseen at the seminar was the

development of some dynamic data sharing standard that can be implemented in several software environments, allowing fast communication between those environments and facilitating software reuse.

Judging by the liveliness of the discussions and the number of joint research projects proposed at the end of the seminar, this meeting between the machine learning and the information visualization communities was more than needed. The flexible format of the Dagstuhl seminars is perfectly adapted to this type of meeting and the only frustration perceivable at the end of the week was that it had indeed reached its end. It was clear that researchers from the two communities were starting to understand each other and were eager to share more thoughts and actually start working on joint projects. This calls for further seminars...