

Information Retrieval Perspective to Interactive Data Visualization

J. Peltonen¹, M. Sandholm¹, and S. Kaski^{1,2}

¹Helsinki Institute for Information Technology HIIT, Department of Information and Computer Science, Aalto University, Finland

²Helsinki Institute for Information Technology HIIT, Department of Computer Science, University of Helsinki, Finland

Email: {jaakko.peltonen, samuel.kaski}@aalto.fi, msandhol@gmail.com

Abstract

Dimensionality reduction for data visualization has recently been formulated as an information retrieval task with a well-defined objective function. The formulation was based on preserving similarity relationships defined by a metric in the input space, and explicitly revealed the need for a tradeoff between avoiding false neighbors and missing neighbors on the low-dimensional display. In the harder case when the metric is not known, the similarity relationships need to come from the user. We formulate interactive visualization as information retrieval under uncertainty about the true similarities, which depend on the user's tacit knowledge and interests in the data. During the interaction the user points out misses and false positives on the display; based on the feedback the metric is gradually learned and the display converges to visualizing similarity relationships that correspond to the tacit knowledge of the user.

Categories and Subject Descriptors (according to ACM CCS): H.5.m [Information interfaces and presentation (e.g., HCI)]: Miscellaneous

1. Introduction

We study methods for interactive data exploration with scatter plots. Traditional simple plots of feature pairs only reveal part of the structure in multivariate data; to show high-dimensional data on a scatter plot, *nonlinear dimensionality reduction* (NLDR) is often applied. What is relevant in data is normally *not known a priori*; it depends on the user and his/her interest. That is why simple visualization using naive assumptions about what data properties are important to show will not work well; *interactive visualization* should be used to let the user give feedback about what is relevant. In this paper we introduce a novel information retrieval based approach to interactive data exploration with scatter plots.

A recent method [VK07b, VPN*10, PK11] formalizes the case where the user is interested in neighborhood relationships between data points. A static (non-interactive) visualization with scatter plots is formalized as a rigorous *information retrieval task* where the user retrieves neighborhood relationships based on the display; the display is optimized to minimize errors between retrieved neighbors and known neighborhoods in the input space. The optimized dis-

play is then a faithful representation of the data in the well-defined sense of yielding few errors in the visual information retrieval. The formalism yields the Neighbor Retrieval Visualizer (NeRV) method which has outperformed several methods [VPN*10]; the methods Stochastic Neighbor Embedding (SNE; [HR02]) and t-Distributed SNE [vdMH08] can also be interpreted as special cases of the formalism.

A general way to encode what aspects of data are relevant to a user is to define the *metric between data*, so that differences between data depend on the relevant aspects. Most NLDR methods require known distances or a known metric; NeRV and related methods use the known metric to compute input space neighborhoods between high-dimensional data. When the metric is not known a priori, a natural approach is to learn it by interaction with the user. We assume the user's interaction is based on an underlying metric that encodes the user's tacit knowledge and interests in the data; we call the input neighborhoods in this tacit metric *true neighborhoods*. Learning the metric and compressing data to the display should be optimized for a unified goal of the interactive system. In this paper we give a solution: **we introduce an interactive visualization method optimized to serve**

the user in the rigorous task of retrieving true neighbors from the scatter plot. We infer the metric iteratively from feedback of the user’s retrieval task, and optimize the display for each iteration, compressing data to the display to serve information retrieval in the inferred metric. In this paper we concentrate on pairwise similarity and dissimilarity feedback. We introduce a mathematical formulation of interactive visualization as *information retrieval under uncertainty of the user preferences*, and our method can be seen as the full interactive extension of the NeRV formalism.

Earlier methods for interactive data exploration with scatter plots include, for example, the Grand Tour [Asi85], and simple approaches where the user explicitly decides what dimensions are selected to be plotted. Recently, systems that try to learn from observation-level interactions how the user thinks data should be arranged [EHM*11, EFN12, BLBC12] have also been proposed. An advantage of our system is that the whole interactive process is optimized for the rigorous user task of neighbor retrieval.

2. New method: information retrieval approach to interactive visualization

Scatter plot visualization of multivariate data is often done by applying nonlinear dimensionality reduction (NLDR; methods include e.g. [TdsL00, BN02, RS00]). Many NLDR methods do not perform well in visualization tasks [VK07a]; they have not been designed to reduce dimensionality beyond the effective dimensionality of the data manifold, and are not good at compressing data onto a low-dimensional display. We first review a formalization of NLDR which has proven successful in static visualization [VPN*10], and we then extend the formalization to interactive visualization.

Let $\{\mathbf{x}_i\}_{i=1}^N$ be a set of input data samples. Let each sample i have an unobserved *true neighborhood* p_i , which is a distribution telling for each neighbor j the probability $p_{j|i}$ that j is chosen as a neighbor to i . The user’s true neighborhoods will be learned from feedback. The goal is to create output coordinates $\{\mathbf{y}_i\}_{i=1}^N$ for the data suitable for visual neighbor retrieval. On the display an *output neighborhood* q_i can be defined around each sample as probabilities $q_{j|i}$, in this paper $q_{j|i} = \frac{\exp(-\|\mathbf{y}_i - \mathbf{y}_j\|^2 / \sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{y}_i - \mathbf{y}_k\|^2 / \sigma_i^2)}$, where $\|\cdot\|^2$ is squared Euclidean distance on the display; $q_{j|i}$ is the probability that an analyst starting from a central point i picks neighbor j for inspection. This simple mathematical form can be replaced by more advanced user models if available.

All properties of high-dimensional data cannot be represented on a low-dimensional scatter plot. Two kinds of errors will happen (Fig. 1, top): *misses* are true neighbors of a point i (high $p_{j|i}$) that are not neighbors on the display (low $q_{j|i}$). *False neighbors* are neighbors on the display (high $q_{j|i}$) that are not true neighbors (low $p_{j|i}$). Misses and false neighbors can have a different cost to the analyst. The display should be optimized to minimize the total cost of errors.

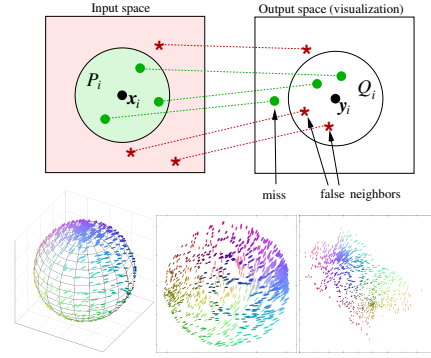


Figure 1: Top: errors in visual information retrieval for query point i . P_i denotes points with high true neighborhood probability, Q_i denotes points with high neighborhood probability on the display. Misses are true neighbors that are not neighbors on the display; false neighbors are neighbors on the display that are not true neighbors. **Bottom:** different tradeoffs between recall and precision (misses and false neighbors) yield different optimal 2D displays. Original 3D data (bottom left) are on a sphere surface; flattening the sphere (bottom center) avoids misses but yields false neighbors from opposite sides, cutting the sphere open (bottom right) avoids false neighbors but yields misses over the cuts. Figure used by permission of [VPN*10].

It has been shown [VPN*10] that the total cost of misses corresponds to the information retrieval measure *recall*, and the total cost of false neighbors corresponds to *precision*. The measures have been generalized to divergences between probabilistic neighborhoods [VPN*10]: the Kullback-Leibler divergence $D(p_i, q_i) = \sum_{j \neq i} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$ is a generalization of recall and $D(q_i, p_i) = \sum_{j \neq i} q_{j|i} \log \frac{q_{j|i}}{p_{j|i}}$ is a generalization of precision. The total information retrieval cost C_{NeRV} of misses and false neighbors is then

$$C_{\text{NeRV}} = \lambda \mathbb{E}_i [D(p_i, q_i)] + (1 - \lambda) \mathbb{E}_i [D(q_i, p_i)] \quad (1)$$

where \mathbb{E}_i denotes expectation over the query points i . The parameter λ in (1) controls the precision-recall tradeoff desired by the analyst: whether misses or false neighbors are more important to avoid. Different tradeoffs yield different optimal low-dimensional displays as shown in Fig. 1 (bottom). In our interactive sessions we emphasize precision (λ near 0) since then intermediate plots are locally well arranged with few false neighbors; this can make it easier to browse data on the display as the analyst is not distracted by false neighbors.

To optimize a scatter plot visualization, (1) must be optimized with respect to the output coordinates \mathbf{y}_i that define the output neighborhoods $q_i = \{q_{j|i}\}$. The previous static visualization approach [VPN*10] can do this only if the true neighborhoods $p_i = \{p_{j|i}\}$ for each data point i are known. In this paper we treat the more difficult case when the true

neighborhoods are unknown; we now extend the approach to unknown true neighborhoods.

2.1. Interactive visualization optimized for information retrieval under uncertainty

Equation (1) can be computed only if the true neighborhoods $p_i = \{p_{ji}\}$ are known. When the true neighborhoods are *unknown*, but evidence of them is available in the form of user feedback, the rigorous approach is to treat the true neighborhoods as missing values, and optimize the expectation of the cost function over the missing values. That is, we *optimize the visualization for information retrieval under uncertainty of the user's preferred similarities*. This is written as

$$\mathbb{E}[C_{\text{NeRV}}] = \mathbb{E}_{\{p_i\}|F} [\lambda \mathbb{E}_i[D(p_i, q_i)] + (1 - \lambda) \mathbb{E}_i[D(q_i, p_i)]] \quad (2)$$

where $\mathbb{E}_{\{p_i\}|F}$ denotes expectation over the possibilities for different true neighborhood distributions $\{p_i\}$; the expectation is over a posterior distribution of the possible neighborhood distributions, given the evidence from feedback F .

Assume the true neighborhoods $p_i = \{p_{ji}\}$ have a simple functional form: they are a function of distances in an *unknown metric* of the original multivariate feature space, so that

$$p_{ji} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_{\mathbf{A}}^2 / \sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|_{\mathbf{A}}^2 / \sigma_i^2)} \quad (3)$$

where $\|\mathbf{x}_i - \mathbf{x}_j\|_{\mathbf{A}}^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{A} (\mathbf{x}_i - \mathbf{x}_j)$ and \mathbf{A} is the metric matrix. Then the expectation over possible true neighborhoods in (2) reduces to an expectation over the possible true metrics, so that

$$\mathbb{E}[C_{\text{NeRV}}] = \mathbb{E}_{\mathbf{A}|F} [\lambda \mathbb{E}_i[D(p_i, q_i)] + (1 - \lambda) \mathbb{E}_i[D(q_i, p_i)]] \quad (4)$$

where the true neighborhoods $p_i = \{p_{ji}\}$ are now functions of the true metric which we denote by its associated matrix \mathbf{A} , and $\mathbb{E}_{\mathbf{A}|F}$ denotes expectation over a posterior distribution of metrics \mathbf{A} given the feedback F .

Equation (4) is an expectation (integral) of C_{NeRV} over the posterior distribution of metrics. In this paper we make a simple fast estimate of the integral. We infer a *variational approximation* $\hat{p}(\mathbf{A}|F)$ to the posterior $p(\mathbf{A}|F)$ as described in Section 2.2. We then approximate the integral by the value of C_{NeRV} at the mean $\mathbf{A}^* = \mathbb{E}_{\hat{p}(\mathbf{A}|F)}[\mathbf{A}]$ of the variational posterior. Since the variational distribution is unimodal and is optimized to contain a large part of the posterior mass, the value of C_{NeRV} at \mathbf{A}^* is a reasonable quick-and-dirty approximation to the integral. We thus write

$$\mathbb{E}[C_{\text{NeRV}}] \approx [\lambda \mathbb{E}_i[D(p_i, q_i)] + (1 - \lambda) \mathbb{E}_i[D(q_i, p_i)]]_{\mathbf{A}=\mathbf{A}^*} \quad (5)$$

where the true neighborhoods p_i are computed by (3) using the mean posterior metric \mathbf{A}^* and the output neighborhoods q_i are computed from display coordinates $\{\mathbf{y}_i\}$ of the data

as defined in Section 2. Equation (5) **measures the performance of a visualization in the information retrieval task of retrieving the true neighbors, corresponding to the analyst's tacit knowledge, from the display**. Equation (5) can be used as an *optimization criterion*, since it is a well-defined function of the display coordinates of the data.

Interactive optimization of the cost (5) performs the following three steps at each iteration. **1.** Infer the approximate posterior mode \mathbf{A}^* of the metric from feedback received so far. **2.** Optimize the visualization for the neighborhoods p_i yielded by the metric \mathbf{A}^* . **3.** Show the new visualization and gather feedback from the analyst. The optimization of the visualization can be done simply by minimizing the cost (5) with respect to each low-dimensional output coordinate \mathbf{y}_i of each data point i ; here we optimize the output coordinates by conjugate gradient descent. The approach has a rigorous interpretation: the display is optimized for minimal expected cost of misses and false neighbors. We call the resulting interactive visualization method the *Interactive Neighbor Retrieval Visualizer* (Interactive NeRV).

2.2. Inference of the metric from feedback

We assume the analyst gives feedback on pairs of points, labeling them similar or dissimilar. We use a Bayesian approach to learn the metric from feedback. The metric is parameterized as $\mathbf{A} = \sum_{d=1}^D \gamma_d \mathbf{v}_d \mathbf{v}_d^\top$ where the \mathbf{v}_d are basis vectors for the data and γ_d are weighting parameters that differentiate the possible metrics. In experiments we use the original basis of the data, therefore we learn weighted Euclidean metrics, which makes analysis of the results easy. Inferring $\hat{p}(\mathbf{A}|F)$ from a set of feedback pairs $F = \{(i, j, f_{ij})\}$ is then done by inferring the variational posterior approximation for the γ_d . The likelihood of a single feedback pair is defined as $p(f_{ij}|\mathbf{x}_i, \mathbf{x}_j, \mathbf{A}, \mu) = (1 + \exp(f_{ij}(\|\mathbf{x}_i - \mathbf{x}_j\|_{\mathbf{A}}^2 - \mu)))^{-1}$, where μ is a threshold parameter and $f_{ij} = 1$ for a similar pair and -1 for a dissimilar pair. Given a Gaussian prior for the weighting parameters γ_d , and the likelihood terms for all feedback pairs as above, we can infer a variational approximation for the posterior of γ_d . Details of the update equations are omitted for brevity. Equivalent Bayesian updates have been used for metric learning without a visualization context [YJS07], but our novel contribution for the metric learning is to integrate the updates as a part of interactive optimization of the information retrieval cost $\mathbb{E}[C_{\text{NeRV}}]$.

3. Experiments

We evaluate our method in three ways: **1.** we evaluate the benefit of utilizing a visualization in finding good feedback pairs, **2.** we test whether the iterative interaction and metric learning help the user in the task of visual retrieval of relevant neighbors, and **3.** we present a small case study with a real user. In experiments 1 & 2, in each iteration 3 pairs of feedback are produced by an artificial mechanism: we com-

pare the current visualization to known true neighborhoods and give the worst misses or false neighbors as feedback.

We use three data sets in the experiments: articles published by researchers of a local research institute, a subset of the DARPA TIMIT phoneme data, and Wine from UCI machine learning repository. Each data set has additional noise features, assumed not to be beneficial for retrieving the true neighborhood relationships corresponding to user interests.

To evaluate our approach we built a simple implementation, where the user is shown a scatter plot and he interacts by picking neighbors and non-neighbors, and can inspect data items by hovering over them with a mouse. Our approach can naturally be integrated in larger systems and can be combined with inspection tools, linked displays, glyphs etc. as in all scatter plot based systems, and with tools to annotate points or regions to ease exploration.

Figure 2 (top) shows, using an oracle user who always picks out the worst miss or false neighbor (with respect to a known true neighborhood), that giving the feedback based on the visualization improves metric learning compared to picking the pair randomly.

Figure 2 (bottom) shows that quality of the visualization improves as the worst misses and false neighbors are pointed out, and that our information retrieval-based visualization approach outperforms traditional multidimensional scaling (MDS) coupled to metric learning on two data sets. Here we measure for each visualization the area under the precision-recall curve, where the true neighbors are defined as the 20 closest neighbors using the ground truth metric. The curves are constructed by varying the number of neighbors retrieved from the visualization between 1 and 100, calculating mean precision and recall for each number of retrieved neighbors. This experiment can be seen as a case of transductive learning: by giving feedback to only a small amount of pairs the overall accuracy of neighborhoods improves.

Figure 3 shows a small-scale user study using scientific articles as the data set. The user’s goal was to arrange the scatter plot in such a way that scientific documents that the user considers similar are close to each other on the screen, by giving pairwise feedback. To help the user browse the points, we displayed the title, year, and authors of the paper in a pop-up display when the user hovered over the corresponding point with the mouse. Additionally, points changed color after the user gave feedback on them. Figure 3 shows that as feedback was given, the metric improved and articles became arranged according to research fields.

4. Conclusions

We introduced a novel interactive visualization method that serves a user in an information retrieval task of finding neighborhood relationships that correspond to the tacit knowledge of the user. The true neighborhoods are encoded

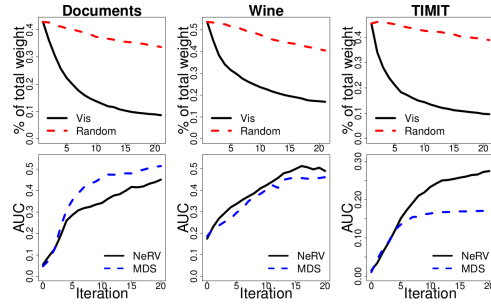


Figure 2: Top: Feedback pairs chosen based on the visualization (“Vis”) improve the metric learning compared to selecting the pair randomly in all data sets. The data sets include added irrelevant features (dimensions) containing only noise, and we measure the portion of weight that the metric assigns to the noise features. Both mechanisms are able to decrease the importance of noise features, our method “Vis” does so faster. **Bottom:** Area under the precision-recall curves using the ground truth metric. Retrieval performance improves when we give feedback about the errors of previous visualizations. Interactive NeRV outperforms the MDS based system on two data sets, but improvement can be seen also with MDS.

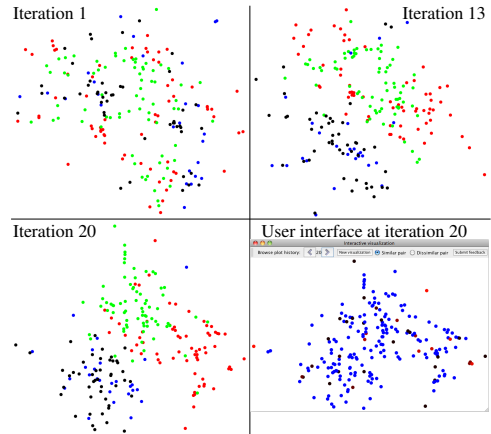


Figure 3: Example visualizations from the user experiment. Colors denote broad fields of research: green=machine learning, red=complexity theory, black=human-computer interaction and blue=social psychology. Starting from the initial visualization (top left), points become more arranged according to the hidden colors as we give more feedback (top right, bottom left). The interface (bottom right) shows data without hidden labels, and feedback point pairs in red shades (brighter red is more recent).

with a metric for the high-dimensional data. The user interacts by pointing out misses and false positives on the display, the metric is inferred from the feedback, and the display is optimized for information retrieval in the inferred metric; the display then iteratively converges to showing similarity relationships relevant to the user. The whole system is rigorously quantifiable and optimizable by performance in the information retrieval task. Our experiments show the interactive visualizer learns metrics better than a simple non-visual mechanism, and shows relevant neighbors better than an alternative multidimensional scaling method coupled to the metric.

5. Acknowledgments

The work was supported by Academy of Finland, decisions 251170 (Finnish CoE in Computational Inference Research COIN), 252845 and 255725. Authors belong to COIN. Documentation data are derived from the following indices: Science Citation Index Expanded, Social Science Citation Index and Arts & Humanities Citation Index, prepared by Thomson Reuters[®], Philadelphia, Pennsylvania, USA, ©Copyright Thomson Reuters[®], 2011, and the Digital Library of the Association of Computing Machinery (ACM).

References

- [Asi85] ASIMOV D.: The grand tour: A tool for viewing multidimensional data. *SIAM Journal on Scientific and Statistical Computing* 6, 1 (1985), 128–143. 2
- [BLBC12] BROWN E. T., LIU J., BRODLEY C., CHANG R.: Dis-function: Learning distance functions interactively. In *IEEE Conference on Visual Analytics Science and Technology (VAST)* (2012), IEEE, pp. 83–92. 2
- [BN02] BELKIN M., NIYOGI P.: Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Advances in Neural Information Processing Systems 14* (Cambridge, MA, 2002), Dietterich T. G., Becker S., Ghahramani Z., (Eds.), MIT Press, pp. 585–591. 2
- [EFN12] ENDERT A., FIAUX P., NORTH C.: Semantic interaction for visual text analytics. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (2012), CHI '12, ACM, pp. 473–482. 2
- [EHM*11] ENDERT A., HAN C., MAITI D., HOUSE L., LEMAN S., NORTH C.: Observation-level interaction with statistical models for visual analytics. In *IEEE Conference on Visual Analytics Science and Technology (VAST)* (2011), IEEE, pp. 121–130. 2
- [HR02] HINTON G., ROWEIS S. T.: Stochastic neighbor embedding. In *Advances in Neural Information Processing Systems 14*, Dietterich T., Becker S., Ghahramani Z., (Eds.), MIT Press, Cambridge, MA, 2002, pp. 833–840. 1
- [PK11] PELTONEN J., KASKI S.: Generative modeling for maximizing precision and recall in information visualization. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics* (2011), Gordon G., Dunson D., Dudik M., (Eds.), vol. 15 of *JMLR W&CP*, JMLR, pp. 597–587. 1
- [RS00] ROWEIS S. T., SAUL L. K.: Nonlinear dimensionality reduction by locally linear embedding. *Science* 290 (2000), 2323–2326. 2
- [TdSL00] TENENBAUM J. B., DE SILVA V., LANGFORD J. C.: A global geometric framework for nonlinear dimensionality reduction. *Science* 290 (2000), 2319–2323. 2
- [vdMH08] VAN DER MAATEN L., HINTON G.: Visualizing data using t-SNE. *Journal of Machine Learning Research* 9 (2008), 2579–2605. 1
- [VK07a] VENNA J., KASKI S.: Comparison of visualization methods for an atlas of gene expression data sets. *Information Visualization* 6 (2007), 139–54. 2
- [VK07b] VENNA J., KASKI S.: Nonlinear dimensionality reduction as information retrieval. In *Proceedings of AISTATS*07, the 11th International Conference on Artificial Intelligence and Statistics (JMLR Workshop and Conference Proceedings Volume 2)* (2007), Meila M., Shen X., (Eds.), pp. 572–579. 1
- [VPN*10] VENNA J., PELTONEN J., NYBO K., AIDOS H., KASKI S.: Information retrieval perspective to nonlinear dimensionality reduction for data visualization. *Journal of Machine Learning Research* 11 (2010), 451–490. 1, 2
- [YJS07] YANG L., JIN R., SUKTHANKAR R.: Bayesian active distance metric learning. In *Proceedings of UAI-07, the Twenty-Third Annual Conference on Uncertainty in Artificial Intelligence* (2007), Parr R., van der Gaag L., (Eds.), AUAI Press, pp. 442–449. 3