

Cluster Sculptor, an interactive visual clustering system

P. Bruneau ^{*}, P. Pinheiro, B. Broeksema, B. Otjacques

Centre de Recherche Public - Gabriel Lippmann, 41 rue du Brill, L-4422 Belvaux, Luxembourg

ARTICLE INFO

Article history:

Received 30 November 2013

Received in revised form

10 September 2014

Accepted 17 September 2014

Available online 11 November 2014

Keywords:

Interactive clustering

Dimensionality reduction

Visual clustering

ABSTRACT

This paper describes Cluster Sculptor, a novel interactive clustering system that allows a user to iteratively update the cluster labels of a data set, and an associated low-dimensional projection. The system is fed by clustering results computed in a high-dimensional space, and uses a two-dimensional (2D) projection, both as support for overlaying the cluster labels, and engaging user interaction. By easily interacting with elements directly in the visualization, the user can inject his or her domain knowledge progressively. Via interactive controls, the distribution of the data in the 2D space can be used to amend the cluster labels. Reciprocally, the 2D projection can be updated so as to emphasize the current clusters. The 2D projection updates follow a smooth physical metaphor that gives insight of the process to the user. Updates can be interrupted any time, for further data inspection, or modifying the input preferences. The interest of the system is demonstrated by detailed experimental scenarios on three real data sets.

© 2014 Elsevier B.V. All rights reserved.

1. Introduction

Clustering algorithms are extensively employed in various domains such as data mining, information retrieval and bio-informatics. They provide means to classify unlabeled multivariate items of various data types in an unsupervised manner. Among other use-cases they are used to find genome-wide expression patterns [1], patterns in trajectories [2,3] and similar documents in a text corpus [4]. In order to exploit the full potential of these algorithms [5], interactive visual representations are required for both analysis and communication purposes.

The high-dimensional spaces real-world data sets often lie in are typically harmful to clustering algorithms. In particular, most well-known clustering algorithms (e.g., k -means [5], spectral clustering [6], or EM for the Gaussian mixture [7, Chapter 9]) rely on the Euclidean distance, or some transform of the latter. Unfortunately, this kind of distance suffers from the curse of dimensionality. As the dimensionality increases, the distribution of pairwise distances is shifted towards high values, while its variance remains almost unchanged (see Fig. 1). Pairwise Euclidean distances thus tend to become indistinguishable when the dimensionality increases [7, Section 1.4]. The use of adaptive [8] or locally sensitive [9] measures may alleviate the problem, the study of which remains central in the machine learning community.

On the other hand, the usage of clustering results as a communication tool is promising, but also affected by the reference to these high-dimensional spaces. The latter are indeed challenging for representation, and user understanding. For effective communication and presentation, the results of clustering algorithms are thus often combined with a Dimensionality Reduction (DR) technique. These techniques [10,11] project a high-dimensional data set to a lower-dimensional space, for visualization in two or three dimensions. The reduced data set, with lower dimensionality, can be visualized using scatterplot techniques. Cluster labels are then overlaid on this low-dimensional projection, using a mapping to glyphs or category colors, as illustrated in Fig. 2.

Such visual representations can be used both for building and adjusting a mental map of the data at hand. They can also be an entry point for finer data inspection using brushing techniques [12]. This need may occur when considering the organization of image collections [13], in the context of multimedia retrieval engines, or when reporting public or medical data. In many cases though, the initial clustering of a data set is deemed to be imperfect with regard to a ground truth, or user expectations. Decision bounds may be unsatisfactory, clusters may be multimodal in the 2D space or exhibit outliers. Worse, data attributes might be badly chosen, noisy, or combined in an inappropriate distance function.

Tackling all these issues at once is certainly not realistic. However, the user should at least be allowed to interact with the clustering and DR results in order to investigate them, and partly cope with them. It is thus possible to let the user manually amend a clustering structure by selecting and labeling points in the projection. However, naive rectangle or lasso selections often

^{*} Corresponding author.

E-mail addresses: bruneau@lippmann.lu (P. Bruneau),
pinheiro@lippmann.lu (P. Pinheiro), broeksem@lippmann.lu (B. Broeksema),
ojacque@lippmann.lu (B. Otjacques).

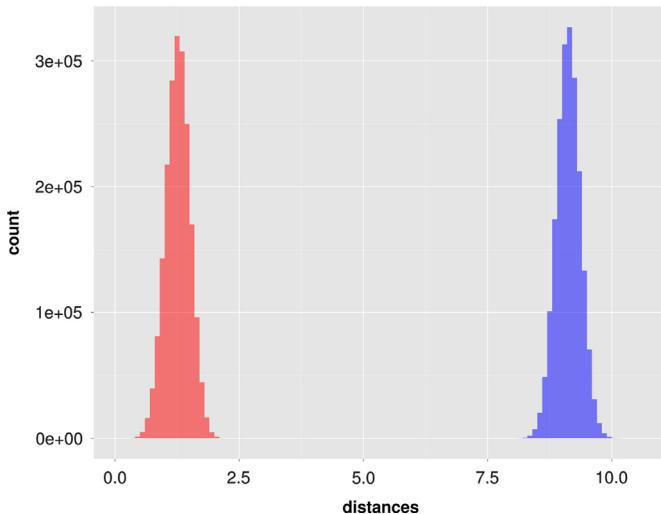


Fig. 1. Distribution of pairwise distances computed from 2000 d -dimensional elements generated uniformly in $[0, 1]$, in red: for $d=10$, in blue: for $d=500$. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)



Fig. 2. t-SNE 2D projection of the COIL-20 image collection. The result from the spectral clustering algorithm is mapped to category colors for the point glyphs. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

result in density gaps within the visual representation of the clusters. This is not consistent with the intuitive meaning of a cluster. Also, when analyzing or communicating clustering results, a user may have preferences in the arrangement of clusters in the visualization space (e.g., highlighting semantic regions). Rearrangement of the clusters' relative positions should thus be supported to some extent. However, clustering and DR techniques usually do not allow such fine tuning.

The purpose of this work is to support the cluster analysis in a visual, semi-automatic way. We largely build upon the t-SNE DR technique [14]. After a review of the related work in Section 2, we recall its use in a batch setting to build an initial 2D projection of the data in Section 3. Then in Section 4, we give a first glance of the Cluster Sculptor workflow and its data inspection facilities. An interactive 2D scatterplot view and legend support the efficient input of user preferences. They are augmented by controls to parametrize and initiate updates of the visualization and the associated clustering. To avoid tedious and error-prone manual grouping and relabeling of data elements, we derive semi-automated

label diffusion techniques, described to a greater extent in Section 5. The user simply has to select few elements (i.e., seeds, in the remainder), which are then used as a basis to interactively refine the clusters. Beyond minimizing the amount of effort the user has to spend, this choice accounts for, and even actively uses, the data distribution in the 2D space. In Section 6, we describe how the underlying dissimilarity matrix is updated behind the scenes to emphasize the clustering structure. Section 7 then shows how the visualization is smoothly and interactively adapted to reflect this update. Detailed experimental scenarios on well known benchmark image collections (COIL-20 [15], MNIST [16]) and an actual biological data set are presented in Section 8. They illustrate tasks that can be conducted with Cluster Sculptor, and demonstrate the interest of the system. After a critical discussion of our proposition in Section 9, we give a summary of our findings, and draw some perspectives for future work in Section 10.

2. Related work

This paper contributes essentially to the interactive and visual clustering state of the art. It assembles ideas taken from the existing literature (i.e., DR technique, clustering, information visualization and label propagation) with crucial contributed parts (e.g., dissimilarity adaptation scheme) in an interactive system.

Therefore, in this section we focus on relating our work to the existing interactive clustering literature. We enrich it with references from several connected domains, such as DR and semi-supervised approaches. The visual and interactive clustering literature covers a variety of work that differs essentially from the pursued objectives.

Seo and Shneiderman propose the Hierarchical Clustering Explorer (HCE) [17]. This system contains various linked views to get oversight and details of hierarchically clustered data, obtained in the context of genetic analysis. They provide means to make cluster comparisons and dynamic query controls to eliminate uninteresting clusters. Their approach is not only focusing on hierarchical clusters, they also seem to assume that the clusters faithfully represent a ground truth as they provide no means to change them.

Turkay et al. [18] also proposed an interactive system to analyze clustering results. With the cluster tendency view and the parallel cluster view, they put an emphasis on the comparative analysis of several clustering results. The parallel cluster view visualizes where data points appear in different clusterings. This view uses the parallel sets technique, which allows us to comparatively find stable structures. As a complement, the cluster tendency view visualizes the similarity matrix of a brushed subset of elements. This allows for validating if the selected data points are likely to be clustered. In our approach we go one step further and provide fine-grained interactive control of the clustering of a selection of data points. Not only do we provide more control over clustering, we also update the projection in order to visually separate the clusters more clearly.

Rinzivillo et al. [19] and Adrienko et al. [3] propose an interactive clustering method for large spatio-temporal data. Initially a density based clustering algorithm is applied to a subset of the data points with a suitable distance function. Next, prototypes are selected for each cluster which, in combination with a cluster distance threshold, form classifiers. Classifiers may be refined by the user by adapting the initial clusters. For example, subclusters can be excluded, turned into new clusters or dissolved among other subclusters. These refinements are supported by visual representations of the clusters and subclusters, which allow the user to interact with both to perform classifier refinement operations. Finally, the obtained classifiers are used to infer the

class of the remaining data points. This approach differs from ours in that it does not have to deal with the placement of clusters, as these are fixed by their geo-spatial coordinates. As a result, not much can be done in this case to make clusters visually more separated.

In the context of document topic modeling, the iVisClustering [4] tool allows us to inspect and interact with textual data and a related latent Dirichlet allocation topic model. Its coordinated views comprise a force-directed layout, based on the similarity of document topic distributions. The topic structure is used to derive cluster summary nodes. The tool is used to identify documents poorly reflected by the current model, or refine vague topics and derive hierarchies of nested topics. The cluster structure is either updated via model parameters tuning or hierarchical refinements. This differs significantly from our approach, where a plain structure is non-parametrically adapted according to user interactions.

Schreck et al. [2] also propose a method for the cluster analysis of trajectory data, but base their approach on Kohonen maps. Traditionally, Kohonen maps are unsupervised. Their initial grid is determined automatically based on, for example, random initialization or principal component analysis of the input data. Schreck et al. propose a user guided approach where some of the grid elements are drawn by the user while the remaining elements are interpolated from the user-specified ones. The reasoning for this approach is very similar to our goal, starting with a sensible initial layout which next can be iteratively refined by the user. Alternatively, we use the unsupervised t-SNE to get an initial map. Also, very similar to what we do is the visualization of the Kohonen map iterative training process. For each step they update the color coding of the cell, allowing the user to visually inspect how the learning process evolves and when it starts to converge. We do something similar by updating the positions of the data points in the map after each t-SNE iteration. The user can interact with the mapping process, and visually estimate when it is stable enough.

ManiMatrix by Kapoor et al. [20] provides an interactive way to change the trade-off in the errors of a classifier. In a visual confusion matrix the user can change the distribution of classification errors based on the requirements for the scenario at hand. Although our approach is in the domain of clustering, it somewhat echoes the approach of Kapoor et al. The ManiMatrix approach starts with initial classifiers which are refined iteratively by the user. Our work presents a similar iterative workflow, but applied to a clustering model.

Contributions in visual clustering also emerged in literature from the DR domain. For example, Broeksema et al. [21] combine dimensionality reduction and clustering in their tool. The applied clustering is based on a Voronoi partitioning of the points in the projection space. The Voronoi cell of a data element is the subset of the projection space such that the data element glyph is closer than any other glyph to all points in the cell (see Fig. 3). Starting at the smallest cell, each cell is merged with its neighbors as long as these neighbors are within a user configured distance. Their approach does not allow for updating the clusters nor the projection to reflect changes in the clustering. Moreover, they focus on categorical data, whereas the scope is set on multi-dimensional numerical data in this paper.

Some interactive DR techniques are also closely related to our work. For instance, Philippeau et al. [22] propose an interactive DR technique for organizing multimedia documents in a 2D visual space. A subset of the documents is placed in the visual space by the user. Based on this placement, a similarity measure is trained which is next used to place the remaining documents in the visual space. A similar approach is taken by Mamani et al. [23], where the differences between default positions of sampled elements and those set by the user feed a feature transform optimization

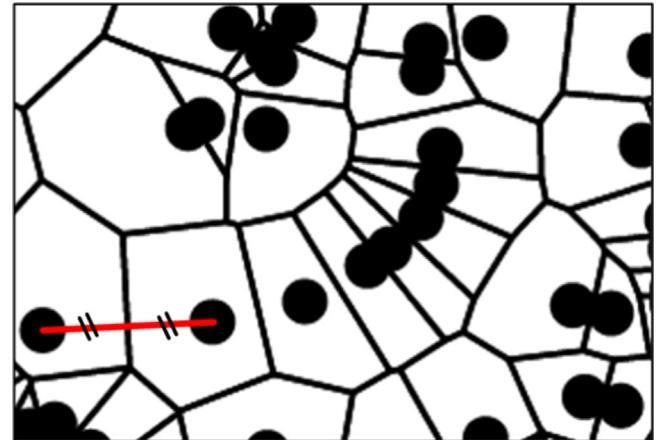


Fig. 3. Voronoi cells examples. The path between two elements is highlighted in red, emphasizing that the respective Voronoi cell border is exactly halfway between them. The glyph in the cell is thus the closest among all glyphs to any points in the cell. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

scheme. Those approaches are different to ours in that the clustering is implicit, i.e., based on placements by the user.

Aupetit emphasizes that all but trivial data sets lie on high-dimensional (HD) manifolds. As such, it is impossible to map such data sets faithfully to a 2D space, as required for rendering scatterplots on-screen. In practice, this causes two types of *projection artifacts* [24]:

- *tears*: two elements close in the HD space are rendered too far in the 2D space,
- *false neighborhoods*: two elements far in the HD space are rendered as neighbors in the 2D space.

Aupetit then illustrates that each DR technique tends to favor a kind of artifacts. For example, PCA [25] is prone to false neighborhoods, whereas CCA [26] is likely to tear manifolds of the HD data. To make these patterns apparent to a viewer, the author suggests to use hovering interactions. Specifically, hovering over a reference element in the 2D space triggers the interactive mapping of the HD distances with respect to other elements. A gray scale coloring of the respective Voronoi cells in the 2D space implements this mapping. This initial work was then extended in ProxiViz [27]. Here, Shepard interpolation [28] is used to smoothen the Voronoi cells, and timer techniques to avoid the flickering that often occurs with false neighbors. This important point pertains to our approach, and we thus included a simplified version of ProxiViz in our system, as shown in Section 4. Along with other tools, such as a data inspector, its use as a support for taking informed decisions before modifying the visualization or the clustering structure is illustrated in Section 8.

Martin et al. [29] propose an approach that has some similarities to our own. In their approach, a user can interactively pose additional constraints on the position of data items in the projected space. These constraints will be taken into account in a next iteration of the DR process. They assume, like we do, that the user can inject knowledge on the similarity of objects. However, they do not concern themselves with providing the users with means to actually cluster objects, which is an explicit goal of our work. Akin to the latter is Dis-Function [30], where user constraints are injected in a feature weight optimization. Full interactive controls are provided, but each iteration requires a full optimization to proceed before the user can see the result of his actions. The authors indicate that a couple of seconds are needed to update the visualization of data sets with barely a hundred

elements, thus not pleading for the scalability of the method's interactivity. Alternatively, we ground upon the physical metaphor underlying t-SNE, to engage the user in following the progressive update consecutive to his or her actions.

Some semi-supervised techniques also use constraints between elements, not necessarily user-specified, to improve a classification function [31]. In this paper, we actually adapt a technique from the semi-supervised literature, the label propagation [32]. We incorporate it as a label diffusion tool. To this respect, our work distinguishes from Dis-Function and the approach by Martin et al. Whereas seeds selected by the user directly affect the distance function and thus the visualization in their work, label diffusion from seeds only modifies the clustering structure in ours. The visualization then uses the clustering to guide its updates more globally.

The present work is roughly a follow-up on the proposition of Bruneau and Otjacques [33]. However, a greater care is now taken about providing consistent user experience, with the use of t-SNE [14]. The technique, overviewed in the next section, allows a smoothly evolving mental map.

3. A summary of the t-SNE projection technique

Let us consider a set of N elements, which is fed as input to Cluster Sculptor. We assume it is defined by d numerical features, with $d > 3$ for HD data sets. Each element is stored as a row of the $N \times d$ matrix \mathbf{X} .

Multidimensional projections aim at representing a HD data set in a visual space (i.e., 2D or 3D) while, to some extent, preserving pairwise distances. In this work, we focus on 2D projections. The goal is thus to estimate the $N \times 2$ matrix \mathbf{Y} of the respective elements in \mathbf{X} . We use the t-SNE [14] DR technique, shown to be resilient to severe tearing and false neighborhood artifacts (see Section 2 for definitions). We motivate this choice by its suitability for the data sets and tasks at hand. In particular, the resulting 2D projections are especially useful to emphasize clusters and the respective low-dimensional local manifolds [34].

Let us define \mathbf{D} and \mathbf{G} as the $N \times N$ matrices of dissimilarities computed from a function of pairs of elements in \mathbf{X} (respectively \mathbf{Y}). \mathbf{D} is thus computed with reference to the original HD space, and \mathbf{G} is its counterpart in the 2D space. These dissimilarities are initialized with normalized Euclidean distances, but any valid dissimilarity (i.e., in the unit interval range) could be used instead.

t-SNE proceeds by defining probability distributions \mathbf{P} and \mathbf{Q} over values in, respectively, \mathbf{D} and \mathbf{G} . \mathbf{Y} is then estimated so that \mathbf{Q} maps \mathbf{P} as closely as possible in the Kullback-Leibler sense [14]. The main feature of t-SNE arises from the use of the heavy-tailed Cauchy distribution for \mathbf{Q} , whereas a Gaussian is used for \mathbf{P} . Beyond leading to a fairly simple gradient optimization scheme, this effectively focuses t-SNE on modeling close neighborhoods in the HD data. As a counterpart, the optimization is relatively insensitive to higher range structures. t-SNE thus distinguishes from many DR methods that try to faithfully model all distance ranges. Though the latter choice is *a priori* preferable, in practice this often leads to crowded visualizations, with all data points packed in the center of the projection [14]. Alternatively, t-SNE encourages a more efficient use of the 2D space. It incidentally decreases the risk of false neighborhood, as shown in an example in Fig. 4.

4. Cluster Sculptor overview

We first give a high-level overview of the approach to interactive clustering in Cluster Sculptor, summarized in Fig. 5. The data

is initially processed by the t-SNE technique exposed in the previous section, which provides a 2D projection of the data to start with. Traditionally, clustering takes place in the HD space, with its results overlaid on the 2D projection (e.g., as category colors) for off-line inspection.

The originality of Cluster Sculptor is to go beyond this static approach, and allow a user to amend the joint 2D projection and clustering using interactive tools. This way, the user can control the state of his or her clustering results, and improve the fit of the visualization to the latter. As exposed further in the remainder, Cluster Sculptor is not a clustering algorithm *per se*. HD clustering results are used as an input, and amended interactively without requiring the full execution of an actual clustering algorithm. An extensive presentation of clustering algorithms is thus outside the scope of the paper. For example, in Section 8, results from *k*-means [5] and spectral clustering [6] are used, but Cluster Sculptor is agnostic of a specific algorithm. It just requires to be fed with a set of labels mapping the data collection under consideration.

Fig. 5 shows four main interactive entry points in the workflow: scope limitation, seed selection, label diffusion, and dissimilarity update. Scope limitation allows the user to limit the effect of his or her further actions to a subset of the data. This is helpful, for example, when a visually homogeneous component is found with multiple cluster labels in it. Next, in the range of influence implied by the scope limitation, seed selection can be used to identify points in the 2D projection from which to start label diffusion. This sequence causes the update of cluster labels. The current cluster labels can then in turn be used to update the dissimilarity matrix \mathbf{D} underlying the t-SNE computations.

We have implemented a web-based prototype that demonstrates the full pipeline of our proposed method. Fig. 6 gives an overview of the Cluster Sculptor interface, and the provided inspection and interaction tools. The front-end is implemented as a one-page AJAX application, and R runs as a server in the back-end. A node.js [35] middleware manages routine calls by the front-end, and all the required data exchanges. R workspace files act as databases, and maintain a consistent state for the application. We give here a first glance of the interface, and the actions that can be performed using it. Details about label diffusion, dissimilarity transforms, or interactive t-SNE updates, are presented afterwards, respectively in Sections 5–7.

The interface is centered on a scatterplot view (Fig. 6e) that displays the $N \times 2$ matrix \mathbf{Y} output by t-SNE steps. When a data set is selected using the data loading controls (Fig. 6a), the data structure needed for t-SNE to run is initialized asynchronously in the back-end. The user is informed of the currently pending operation. The user may then run or stop the t-SNE process on demand. A label vector, selected from the available clustering results in the data loading controls, triggers the overlay of categorical colors on the current scatterplot. This set of colors is statically defined following the recommendations of Harrower and Brewer [36]. The controls also feature several convenient overloads, such as cached t-SNE results, or updated dissimilarity matrices.

The controls from Fig. 6b influence the appearance and behavior of the scatterplot view content. The user may switch any time between a classical scatterplot view, as shown in Fig. 6e, or a Voronoi diagram view (e.g., Figs. 7a and 20b). Both visual primitives support two hovering interactions. The default behavior is to display a summary of the hovered point in the data inspector (Fig. 7b). The pairwise distances in the HD space relative to the hovered point (Fig. 7a) can also be interactively mapped if the associated control is checked. This interaction is largely inspired by the ProxiViz tool [27], and is a valuable asset to assess the faithfulness of the visualization. Accounting for our context, we incorporated slight modifications to the tool. The currently

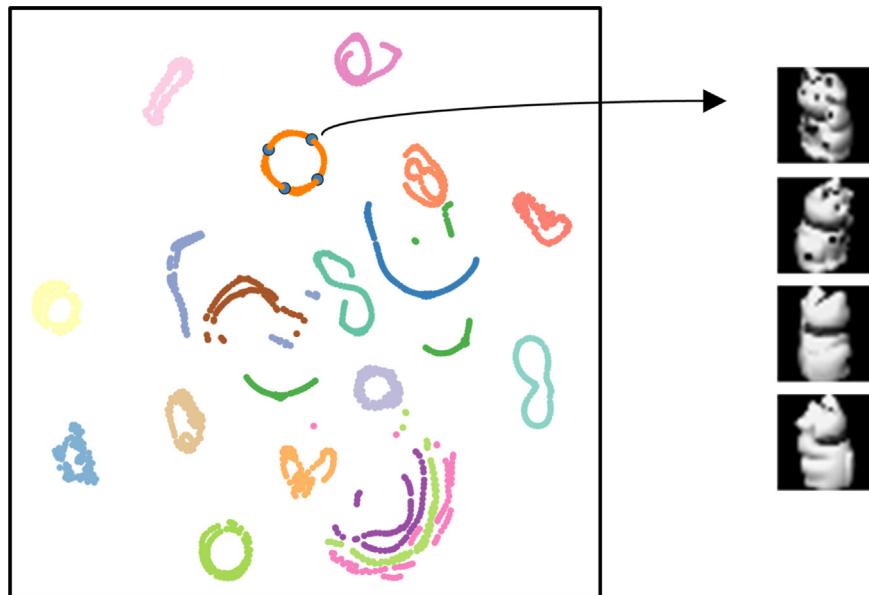


Fig. 4. Left: t-SNE 2D projection of the COIL-20 image collection (see Section 8 for a description). Ground truth classes are mapped to the point glyphs as categorical colors. Right: The highlighted glyphs illustrate the ability of t-SNE to recover the manifold implied by the rotation of the class object. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

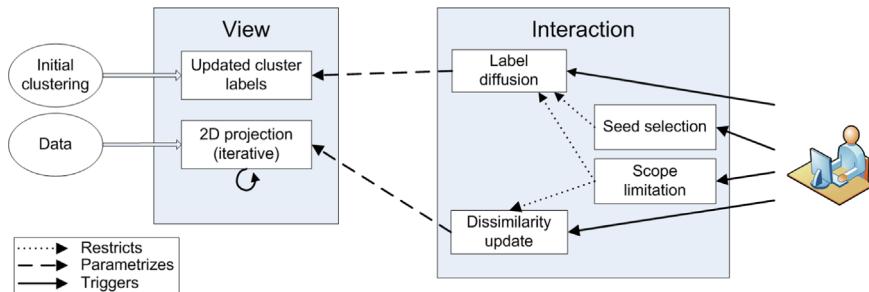


Fig. 5. High-level overview of the Cluster Sculptor approach.

hovered point or cell is emphasized with its respective cluster categorical color (instead of white in ProxiViz). The categorical legend is dynamically substituted with a scale on the distance range when the mouse is inside the scatterplot view (see Fig. 7a). We also added a recall of the cluster labels in the Voronoi diagram view, by overlaying stroke-colored point glyphs on the cells.

The behavior of the t-SNE updates can be controlled via the scatterplot update controls (Fig. 6c). A check-box controls the restriction of t-SNE updates to the current scope; this notion is explained to a greater extent below.

The label diffusion and dissimilarity transform controls are the prime tools for amending the currently displayed clustering and visualization (Fig. 6d). The label diffusion techniques use the distribution of the data in the 2D space to reshape the current set of labels. The dissimilarity transform facilities mirror them, by triggering the update of the HD-related dissimilarity matrix \mathbf{D} underlying the t-SNE computation, based on the current set of labels. The outcome of such transforms can be immediately visualized using the ProxiViz tool, and affect subsequent interactive t-SNE steps.

In addition to hovering, element glyphs can be clicked to define label diffusion seeds. The currently selected seeds are listed (see Fig. 8), and can be interactively edited, for example to define multiple seeds for a single label value. As already mentioned, a legend next to the scatterplot view lists the mapping between glyph colors and cluster labels. Labels can be interactively edited, e.g., to give a semantic value to a cluster, or to trivially merge

clusters when typing an already mapped label. The color patches are also interactive: clicking them highlights the associated cluster in the visualization. Several clusters may be selected simultaneously, and the selection is persistent beyond modifications of the label vector. This can be useful to restrict the comparative analysis of several clustering results (see Fig. 9).

Beyond implementing a simple visual filter of the data according to the cluster labels, the selection mechanism also implicitly defines a scope. In this paper, a scope is a restriction of the elements according to the selected clusters. When a scope is active, operations such as label diffusions, dissimilarity transforms, or even t-SNE updates (if the control in Fig. 6c is checked) occur on this restriction. Elements outside this scope are then by no means affected. With this mechanism, a user can work iteratively on parts of the projection, without excessive discontinuities with respect to his mental map. This limits the cognitive burden and the visual clutter. Behind the scenes, the computational burden of diffusion and transform operations is also decreased.

5. Label diffusion from selected seeds

As explained in the previous section, reference elements, i.e., *seeds*, can be selected directly in the scatterplot. They can be seen as tentative labels, associated to a single element at the moment of the selection. Their purpose is to be later diffused in a chosen scope, according to a specific mechanism. We see them as a

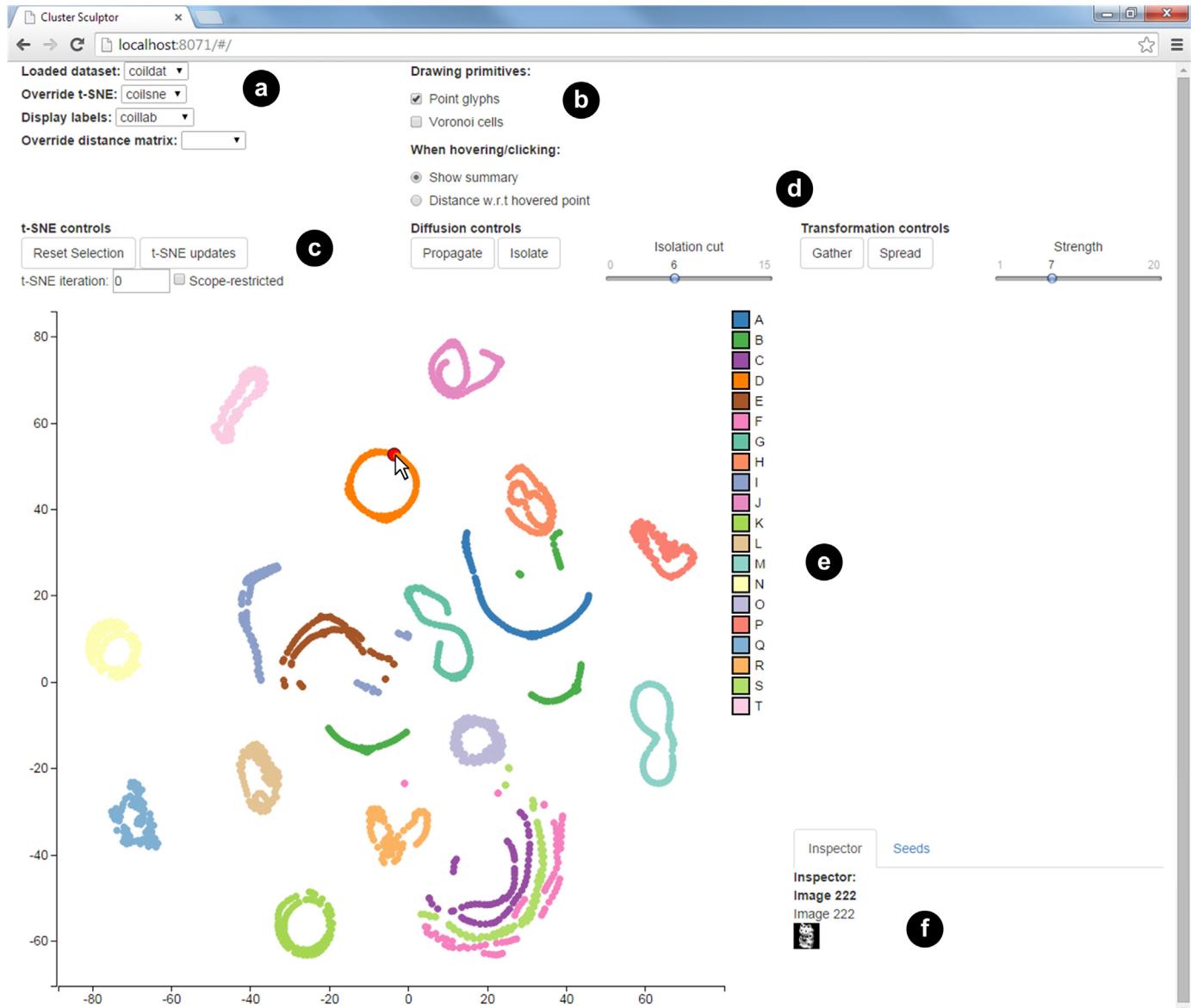


Fig. 6. Overview of the Cluster Sculptor interface. (a): Controls to load data from the back-end. (b): Controls of the element aspect in the scatterplot, and the behavior when hovering/clicking. (c): Scatterplot update controls. (d): Label diffusion and dissimilarity transform controls. (e): The scatterplot, and the associated interactive legend. (f): The data inspector, stacked to the interactive list of currently selected seeds.

semi-automatic alternative to the classical rectangular and lasso selectors, that uses the 2D distribution of elements to update the clustering structure. We implemented two diffusion mechanism, that are of complementary use, as shown experimentally in Section 8. The probabilistic *label propagation scheme*, described in Section 5.1, is inspired by the semi-supervised learning literature. It performs a comprehensive diffusion of the seeds in the current scope. The *isolation scheme* (Section 5.2) uses the Minimum Spanning Tree of the current scope. A breadth-first exploration in cuts of the tree induces a limited diffusion of the seeds.

5.1. Label propagation

This operation is an interactive adaptation of the label propagation technique [32], taken from the semi-supervised learning literature. In this section, for mathematical convenience the seeds selected in the current scope are assumed to take values in $1 \dots C$. This set of seeds is described here as the labeled set, whereas the

remainder of the scope is the unlabeled set. The goal of the operator is to propagate the known labels exhaustively.

Consistently to the definitions in Section 3, let us define \mathbf{Y}^s as the restriction of \mathbf{Y} to the current scope. \mathbf{Y}_L^s (respectively \mathbf{Y}_U^s) is then the labeled (respectively unlabeled) subset of l (respectively u) data elements. For further convenience, elements in \mathbf{Y}^s are permuted so that:

$$\mathbf{Y}^s = \begin{bmatrix} \mathbf{Y}_L^s \\ \mathbf{Y}_U^s \end{bmatrix} \quad (1)$$

Each row in \mathbf{Y}^s has a respective counterpart in \mathbf{Z}^s , the probabilistic labels for the elements in the scope. The C columns of \mathbf{Z}^s mirror the C seed values, with \mathbf{Z}_{nc}^s the probability of element n having label c , and $\sum_{c=1}^C \mathbf{Z}_{nc}^s = 1$. Values in \mathbf{Z}_L^s are thus set to binary values, so as to reflect the seeds set by the user (see Fig. 8 for an example).

The propagation scheme metaphorically lets labels *jump* from element to element. It follows the intuitive idea that similar

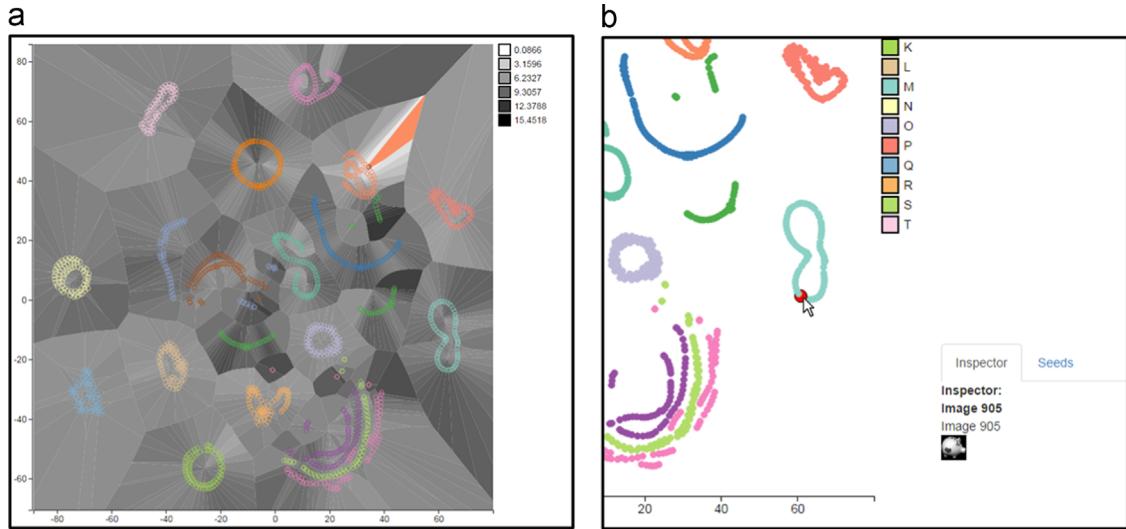


Fig. 7. (a) Voronoi diagram view of the COIL-20 2D projection. The ProxiViz tool is activated, and overlays distances in the HD space w.r.t the hovered point. In this context, the legend shows the mapping of distances in the HD space to the gray scale. (b) Classical scatterplot view, supported by the data inspection tool. The hovered point is highlighted, and the data inspection panel is refreshed interactively.

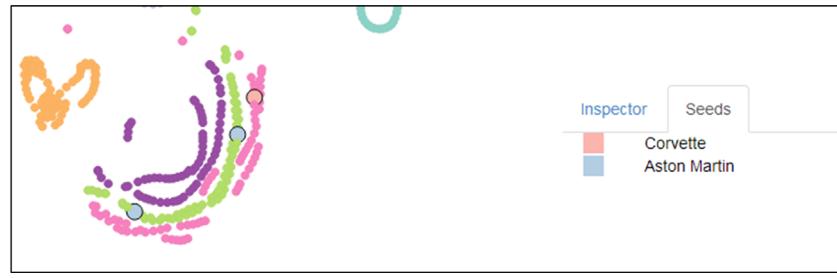


Fig. 8. Panel listing the currently selected seeds. The labels can be interactively edited.

elements are likely to have similar probabilistic labels. Instead of elements described in a vector space, the propagation algorithm thus uses similarity values, ranging in [0, 1]. Such values can be obtained by an unnormalized Gaussian applied on pairwise distances between rows of \mathbf{Y}^s , for example. We note \mathbf{S} the matrix in which $S_{nn'}$ is the similarity between rows \mathbf{Y}_n^s and $\mathbf{Y}_{n'}^s$. We define a probabilistic transition matrix as

$$\mathbf{T}_{nn'} = P(n' \rightarrow n) = \frac{\mathbf{S}_{nn'}}{\sum_{m=1}^{l+u} \mathbf{S}_{mn'}}, \quad (2)$$

with $\mathbf{T}_{nn'}$ being the probability of jumping from element n' to element n . Mirroring the permutation defined by Equation (1), \mathbf{T} has the following block structure:

$$\mathbf{T} = \begin{bmatrix} \mathbf{T}_{ll} & \mathbf{T}_{lu} \\ \mathbf{T}_{ul} & \mathbf{T}_{uu} \end{bmatrix} \quad (3)$$

To avoid notation clutter in later steps, \mathbf{T} is assumed to have its rows normalized to a unit sum. The propagation then proceeds by iterating $\mathbf{Z}^s = \mathbf{T}\mathbf{Z}^s$ until \mathbf{Z}^s converges. In [32], the authors showed that this algorithm converges to a unique fixed point, and that \mathbf{Z}_U^s can be initialized arbitrarily without influence on this fixed point. Specifically, the converged solution is shown to be

$$\mathbf{Z}_U^s = (\mathbf{I} - \mathbf{T}_{uu})^{-1} \mathbf{T}_{ul} \mathbf{Z}_L^s \quad (4)$$

An example application of this operator is shown in Fig. 10. Applying Eq. (4), the operator thus diffuses the C seeds

exhaustively to the data set \mathbf{Y}^s , following the topological information provided in \mathbf{S} .

5.2. Isolating manifolds using the minimum spanning tree

The propagation operator described in Section 5 diffuses exhaustively the seeds to the scope. Cluster Sculptor also features a more exclusive tool, that isolates local manifolds in the scope. It also proceeds from the seeds. Similar to the label propagation scheme, it diffuses them according to the vicinity between elements in the visualization, as implied by the pairwise distances in the 2D space.

Let us consider the complete graph over the elements in the scope. Borrowing notations from Section 5.1, we define $E_{nn'}$ the edge between elements \mathbf{Y}_n^s and $\mathbf{Y}_{n'}^s$. This edge is weighted with $w_{nn'}$, the Euclidean distance between \mathbf{Y}_n^s and $\mathbf{Y}_{n'}^s$. The Minimum Spanning Tree (MST) is then defined as the subgraph that connects all elements with minimal summed weight (see Fig. 11a). It is usually obtained using Kruskal's algorithm [37], or Prim's algorithm [38]. We also define a cut of this graph at w_{cut} as its restriction such that

$$E_{nn'} \in \text{cut} \leftrightarrow w_{nn'} \leq w_{\text{cut}} \quad (5)$$

Prior to triggering the *isolate* control in the interface, the user may adjust w_{cut} using the *isolation cut* control (Fig. 6d). Updating the cut of a MST is fast, and can be performed interactively (see Section 9 for notes on MST-related computational complexity).

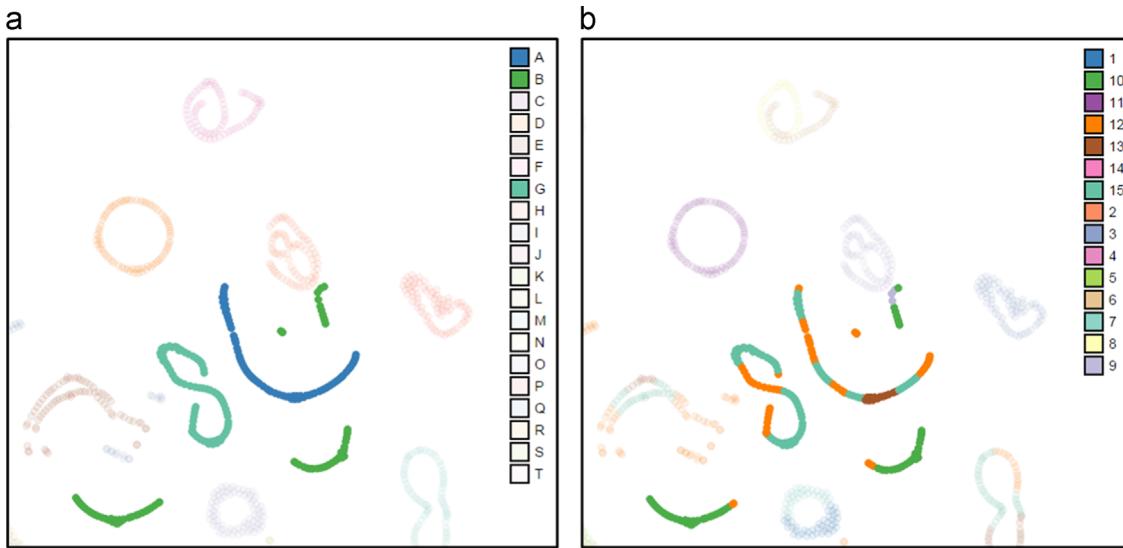


Fig. 9. (a) Example of scope restriction. The selected clusters are highlighted by lowering the alpha channels of non-selected ones in the scatterplot and in the legend. (b) When updating the label vector, the current scope is maintained. The user may reset the selection when this specific scope is not needed any more.

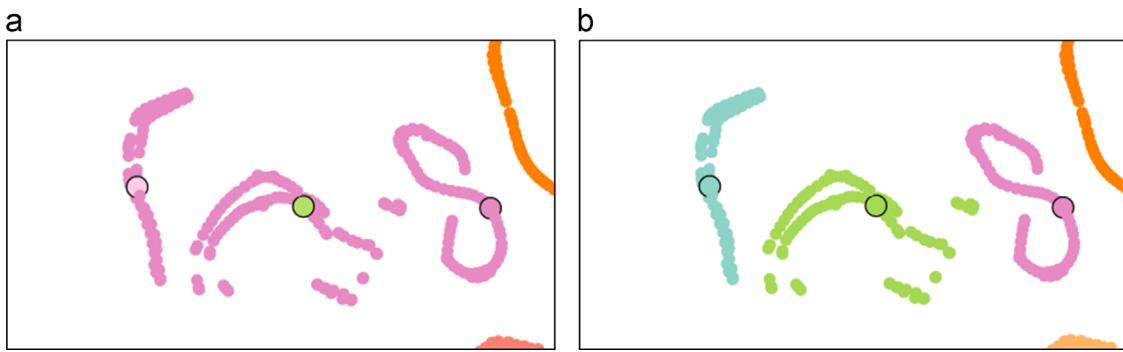


Fig. 10. Example of propagation of 3 seeds. A subsample from the t-SNE projection of COIL-20 is used.

A set of connected components can then be extracted from the MST cut using simple breadth-first explorations. Depending on the set parametrization, triggering the *isolate* control has a variable effect:

- if no seed is selected, root nodes for the breadth-first search are chosen randomly in the scope as long as the scope has not been completely processed. All extracted connected components become new clusters, replacing those that defined the scope,
- if one or more seeds are selected, new clusters are formed with the connected components extracted when using the respective seeds as root nodes of the search. Depending on user preferences as set in the controls, the remaining elements keep their label as prior to scope selection, or are regrouped in a new cluster.

The result of an isolation cut is shown in Fig. 11. After their update, the current cluster labels serve as cues for updating the dissimilarity matrix underlying the visualization. This mechanism is presented in the next section.

6. Dissimilarity updates

A user may wish to update the 2D projection, to reflect the clustering information as well as possible. This may happen either on the account of the clustering results initially loaded, or after

updating the cluster labels using tools described in the previous sections. Multiple reasons can be invoked:

- the clusters have too shallow borders,
- clusters may be multimodal in the 2D space,
- clusters may have unsatisfactory neighborhoods, irrelevant to user knowledge.

Overall, the user could either want to inject knowledge absent from the raw features used to compute the dissimilarities underlying the 2D projection, or emphasize cluster boundaries. In this paper, we choose to translate these preferences in updates to the dissimilarity matrix \mathbf{D} used by t-SNE. This DR technique is exactly about mapping the distribution of values in \mathbf{D} to \mathbf{Y} (see Section 3). Therefore, modifying the dissimilarity matrix is expected to be progressively reflected in the 2D projection as a side effect of t-SNE updates.

In this section, we focus on the dissimilarity modification process. We identify two ways a user would want to update the distribution of the clusters in the visualization space:

1. *gathering* them, i.e., bringing separate clusters in the current scope closer to each other.
2. *spreading* them, i.e., separate clusters in the scope that are close to each other.

Let us consider the complete weighted graph implicitly defined by the matrix \mathbf{D} . We choose to restrict dissimilarity updates to the

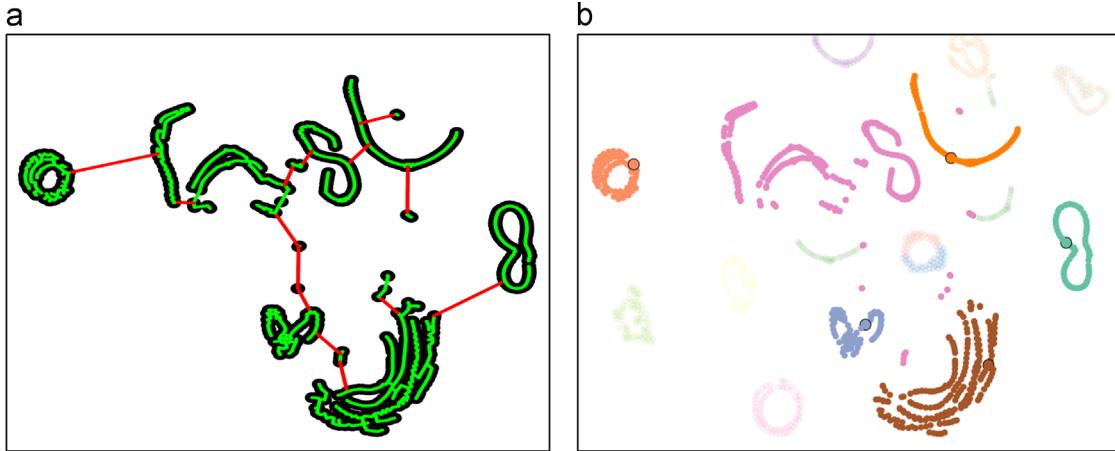


Fig. 11. (a) MST of a subset of a COIL-20 t-SNE projection. The cut for $w_{cut} = 5$ is shown in red. (b) For the highlighted set of seeds, the resulting new clusters are shown. According to the user parametrization in this case, connected components with no associated seed are gathered in a cluster. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

bipartite subgraphs between clusters in the current scope. This choice preserves the topology of components in the visualization. Additionally, the efficient use of the 2D space by t-SNE prevents elements from excessively packing together.

In Cluster Sculptor, both operations are implemented as different parameterizations of the cumulative beta distribution function, $P_{\text{beta}(\alpha,\beta)}$. Considering a scope with two label values l_1 and l_2 , this function is applied on the bipartite graph between l_1 and l_2 . Edges of this graph link elements with label l_1 to elements labeled with l_2 . Generalizing this principle to any number of labels, the associated restriction of \mathbf{D} is filtered according to $P_{\text{beta}(\alpha,\beta)}$:

$$\mathbf{D}_{nn'}^{\text{new}} \leftarrow P_{\text{beta}(\alpha,\beta)}(\mathbf{D}_{nn'}), \quad (6)$$

for n, n' referring to an edge in the bipartite graph. To preserve the structure of local manifolds initially learned by t-SNE, we enforce monotonous updates of the dissimilarities. In other words, we require $\mathbf{D}_{nn'} > \mathbf{D}_{mm'} \rightarrow \mathbf{D}_{nn'}^{\text{new}} > \mathbf{D}_{mm'}^{\text{new}}$. This is done by constraining either α or β to be 1. The resulting family of functions is shown in Fig. 12. This family of functions bijectively maps the unit domain to itself. While maintaining the order of dissimilarities, this guarantees valid dissimilarity values after the filter application. As Fig. 12 shows, curves for $\alpha > 1$ (resp. $\beta > 1$) tend to reduce (resp. increase) the dissimilarity between elements, thus effectively *gathering* (resp. *spreading*) them.

The impact of the operation can be more clearly interpreted in Fig. 13. This figure is showing the absolute increase or decrease of dissimilarity caused by applying P_{beta} for a range of α and β values. Fig. 13a highlights that enforcing valid similarities causes the functions to be bounded from above and below. Then Fig. 13b displays the *closeness to the bound* of the dissimilarities after the transform, i.e., the relative influence of the operation. As gathering and spreading functions are mirroring, these closeness profiles are similar in both cases. From this plot, it is clear that increasing α or β tightens a wider range of similarities to the minimal or maximal valid dissimilarity. The choice of the P_{beta} family of functions is motivated by the behavior at the vicinity of $x=0$. Trying to spread elements that are extremely similar can be disruptive for the visualization. As shown in Fig. 13b, P_{beta} smooths this problem, and ensures the manifolds underlying \mathbf{D} are preserved to some extent. α or β can then be seen as the *sharpness* of this smoothing.

In Figs. 12 and 13, dissimilarities are assumed to range in the unit interval, which may be overly restrictive in practice. Considering an arbitrary domain $[a, b]$, $P_{\text{beta}(\alpha,\beta)}$ can be rescaled

according to the following equation:

$$P_{\text{beta}(\alpha,\beta)}^{[a,b]}(x) = (b-a)P_{\text{beta}(\alpha,\beta)}^{[0,1]}\left(\frac{x-a}{b-a}\right) + a \quad (7)$$

This operation preserves the required properties, i.e., it bijectively ranges in $[a, b]$, and preserves the order of dissimilarities, with minimal and maximal bounds respectively at a and b . In practice, we use this rescaling to ensure two desirable properties:

- the modified dissimilarity should not exceed the current maximal dissimilarity in \mathbf{D} ,
- when gathering two clusters, the decreased dissimilarity should account for the internal cohesion of the clusters.

The first property amounts to set b statically for a given data set. When spreading clusters, a plays a marginal role, and can be left to 0. When gathering clusters, we choose to set a to some quantile of the distribution of dissimilarities internal to the clusters (i.e., outside the bipartite graph). In the experiments described in Section 8, we use the 5% quantile.

The next section describes how we adapt the t-SNE algorithm to discontinuous changes in the matrix \mathbf{D} such as described above, and smoothly render this change to the user in an animated fashion.

7. Updating the t-SNE

Being a gradient-based method, t-SNE is quite expensive to compute afresh [34]. For instance, it took approximately 11 min, on an 8 core machine, to compute 1000 gradient steps for the COIL-20 data set (1440 elements, 30 principal PCA features). Assuming a perturbation of the dissimilarity matrix \mathbf{D} such as described in Section 6, computing a t-SNE projection from the ground up at each user interaction is clearly not acceptable.

The available R implementation [39] is a classical batch algorithm. We adapted it to support independent step executions, and maintain an internal state for the algorithm on the R server, as our interactive scheme requires. Convenience accessors are also implemented to support updates of any part of the internal state. We can for example modify its parametrization, update the dissimilarity matrix (see Section 6), and overload the current 2D projection (see Fig. 6a). Let us note that updating \mathbf{D} triggers the recomputation of the \mathbf{P} matrix that underlies the projection.

t-SNE can be qualified as an *anytime* method, with each iteration being a smooth update of the preceding 2D layout, following a spring-based physical metaphor [14]. Consequently, discontinuous updates of \mathbf{D} are supported by continuously updating the 2D position matrix \mathbf{Y} towards its new convergence point, ensuring the visual stability of the projection. Furthermore, the convergence of gradient-based methods such as t-SNE is typically difficult to assess with automatic means. A user is actually much more qualified to assess visually when the projection is stable enough, and can take this decision using the t-SNE controls.

Spring-based layouts are prone to get stuck in local optima [13]. Worse, a complete graph is involved in t-SNE, which causes a high effort against any live update of \mathbf{D} . In the context of the batch t-SNE, this problem is handled by a simulated annealing phase [14]. Though efficient in the latter case, doing so when live updating \mathbf{D} would be visually disruptive, and thus inadequate.

Instead, we allow the restriction of the gradient steps to the current scope. The gradient sums are then performed only over elements in the selected scope, ignoring the rest of the data set. This option can be activated on user demand (see t-SNE update controls in Fig. 6c). As experimentally shown in Section 8, this

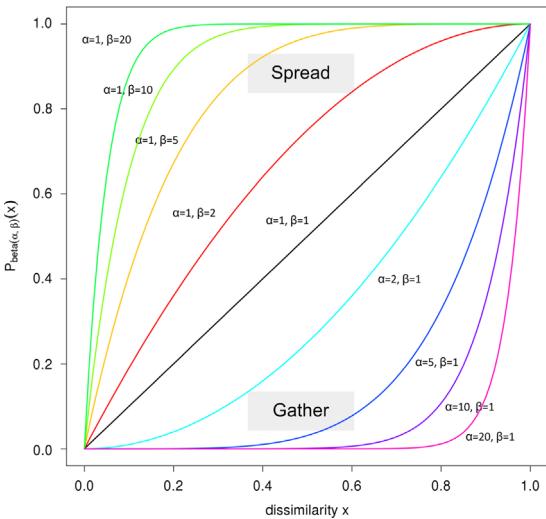
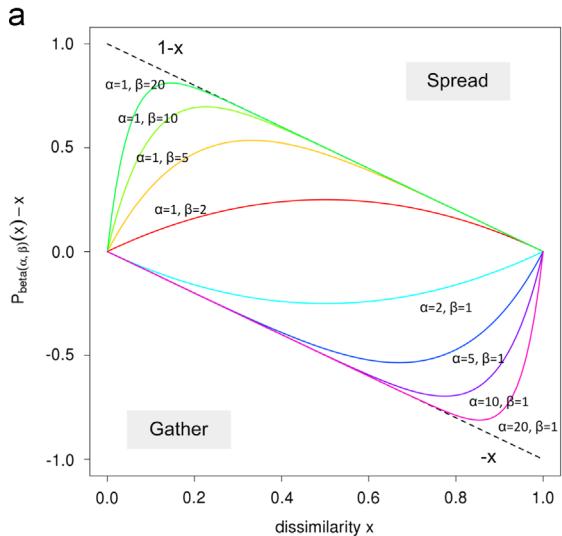


Fig. 12. The family of cumulative beta distribution functions with either α or β set to 1. The identity (both parameters set to 1) is indicated as a reference.



option is effective at quickly updating the position of clusters in the scope. However, their new positions may be conflicting with others outside the scope. To solve this, the user has to find a satisfactory configuration by alternating scope-restricted and classical update sequences. If necessary, this might be supplemented by including the conflicting cluster in the scope, and using an additional spread operation.

Via detailed experimental scenarios on real data, the next section shows how the assembly of tools described above can be used to amend an initial 2D projection, and associated clustering results. We specially put an emphasis on how the tools are used to harness the differential information carried respectively by the 2D projection and the HD clustering results.

8. Experimental scenarios

This section describes three scenarios of a user that iteratively amends initial 2D projections and their associated clustering using Cluster Sculptor. Rather than focusing on low-granularity tasks, we show on a higher level how a user can use the tools. We demonstrate how to get some insight of the data, and how to update the visualization and clustering structure according to his or her findings.

We use two image collections: COIL-20 and MNIST. COIL-20 contains 1440 images, each having $32 \times 32 = 1024$ pixels. The images are photographs of 20 objects rotated by all possible angles modulo 5° . A class is thus made of the 72 images of a given object. The MNIST collection is made of 60,000 images, with each $28 \times 28 = 784$ pixels. The images are variants of handwritten digits, thus forming 10 even classes.

Additionally, we used a biological data set, describing 31,483 bacterial DNA fragments over 8 numerical features. These fragments were sampled from digesters, where biodegradable waste is stored, and slowly degraded by bacteria, to eventually generate biogas. This data set was handed by biologists, in the context of a beginning collaborative exploratory analysis of the data. We further refer to this data set as *biogas*.

MNIST and biogas are too big for the current Cluster Sculptor implementation (see Section 9 for a discussion on complexity issues). Therefore, we sampled randomly respectively 1500 and 3000 elements in these data sets. After removing dimensions that carry no information in each collection (i.e., zero variance), SVD is

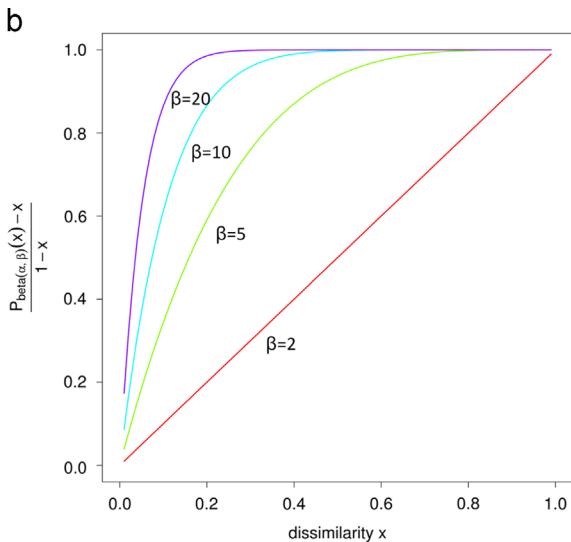


Fig. 13. (a) Plot of $P_{\text{beta}(\alpha, \beta)}(x) - x$ with respect to x . $1 - x$ and $-x$ bound this family of functions from above and below. These bounds ensure transformed similarities remain valid. (b) Closeness to the bound of the spreading functions. Curves for gathering functions mirror them, i.e., decrease from 1 to 0.

applied on these. This removes effects caused by potential correlations among variables. Van der Maaten and Hinton prescribe to retain the coordinates on the 30 principal axes as element representatives [14]. For COIL-20, this amounts to retaining 88% of the data set variance, as indicated by the singular values. For MNIST, 30 principal axes summarize less than 80% of the variance. We thus extended the size of the axes set to 50, that explain 83% of the variance. biogas has only 8 variables, so all axes can be retained.

The next subsections present interaction scenarios specific to each data set. In the scenarios, we refer to clusters with integer IDs, as we have no prior information on their meaning. If desirable, these default IDs can directly be edited in the interactive legend (see Section 4).

8.1. COIL-20

Selecting COIL-20 in the data loading controls (see Fig. 6a) triggers the initialization of the t-SNE internal state. The initial projection is eventually displayed (see Fig. 14a). The user starts the t-SNE iterations by clicking *t-SNE updates* button. Iterations continue to happen until the *t-SNE stop* button is clicked. The result after 1000 iterations is shown in Fig. 14b.

The user then selects a label vector obtained with the spectral clustering algorithm [6] on the HD representation of the data set. At the time of the clustering process, the user had no clue about an adequate number of clusters. As an initial guess the algorithm was parametrized with 15 clusters (Fig. 14c).

The clustering algorithm visibly captured some clear structure in the data set, e.g., Cluster 14 highlighted in Fig. 14c. However, it also fails to recover clear patterns discovered by the t-SNE projection. Maybe resulting from a bad parametrization, some clusters aggregate several unrelated groups. For example, see Clusters 12, 13 and 15 in Fig. 15a. The adjusted Rand index [40] of the clustering w.r.t. the ground truth labels, of 33%, reflects these observations to some extent. Next, the user uses the tool to improve this initial guess.

Before modifying the 2D projection, the user wants to quickly reorganize Clusters 12, 13 and 15 using the label diffusion tools. He or she defines a scope for this operation, by clicking the related color patches in the legend. He or she then selects 5 seeds, purposely to isolate the 5 manifolds these seeds lie on (Fig. 15b). The remainder of the scope is to be regrouped in Cluster 18. After adjusting the *isolation cut* slider, and observing interactively the resulting isolations, the users retains 5 for this parameter. This is leading to the result shown in Fig. 15b. It is worth noting that the short sequence of actions taken until now (a dozen of clicks and few slider interactions) led the adjusted Rand index of the clustering to 69%.

The user then notices Cluster 10, with 3 distinct manifolds in the current clustering. To restrict his or her attention, the user defines a scope on Clusters 9, 10 and 18. Elements of the 3 components are first hovered over. The data inspector is used to confirm that they indeed are related to the same ground truth object (see Fig. 16a). The user then activates the Voronoi diagram view, along with the ProxiViz tool. In the scope restriction, he or she searches potential tears that would suggest a re-unification of the 3 components of Cluster 10. The lighter the gray shade, the closer the mapped element in the respective Voronoi cell is to the reference element in the HD space. Using Cluster 10 and the neighboring ones as a cue for an inspection with the ProxiViz tool, manifold tear artifacts are clearly identified (Fig. 16b and c). The user decides to amend the 2D projection, so that Cluster 10 is displayed as a single manifold.

First, as some confusion errors seem to affect Cluster 10 according to the 2D projection, the user marks each of its components with a seed, whereas one is sufficient for Cluster 9.

The user isolates Clusters 9 and 10 after having adjusted the isolation cut. The remainder of the scope is again assigned to Cluster 18 (Fig. 17a).

Setting the sole Cluster 10 as the scope, the user tries to run restricted t-SNE steps to see if the distinct components can be reunited this way. Two of the three components merge quickly, after approximately 50 iterations. However, the algorithm stabilizes to a configuration with still two distinct components. One of those is in conflicting positions with another cluster (Fig. 17b). Striving for the reunion of the two components, the user defines a seed on each component, and diffuses them with the propagation tool. A gathering transform is then applied to this scope, with moderate sharpness (3), to encourage the merger of the temporary clusters as a single manifold. Subsequent t-SNE steps on this scope indeed quickly lead to a reunion, with its internal structure reflecting the underlying data topology. The temporary clusters are then trivially merged, and restored as Cluster 10, using simple edits in the legend (see Section 4). Few switches between scope-restricted and unrestricted t-SNE steps lead to a proper cluster separation (Fig. 17c).

The user then notices Cluster 13, standing out from the others with its distinct shape (Fig. 18a). It actually regroups the images of three cars. As car images differ more by their pose than by the patterns overlaid on them, t-SNE has laid them in close manifolds.

Distinguishing the three cars happens to be difficult for clustering algorithms, leading to a single cluster. This layout does not conform to user preferences, rather expecting a single car in each cluster. The user thus decides to modify the clustering and 2D projection to have one of the cars standing out from the others. The user sets three seeds and tests several isolation cuts with slider adjustments. This way he is able to derive a close to correct cut of one of the cars in the original Cluster 13 (see Fig. 18b). He or she then applies a spread dissimilarity transform to make the separation clearer. Interestingly, as a side effect, a part of Cluster 18 referring to the same car as the new Cluster 19, but initially lying out of Cluster 13, is automatically connected to the new cluster (see Fig. 18c). Label 19 can then be quickly diffused. This is done by first isolating a seed on the part of Cluster 18 highlighted in Fig. 18c. A trivial merger is then performed by setting 19 as the legend ID of the freshly isolated component.

8.2. MNIST handwritten digits

Distinctly from the COIL-20 scenario, the user loads a cached t-SNE projection of the MNIST data set, that has been computed off-line on the server. The spectral clustering algorithm has also been processed off-line using the HD data, and the result is overlaid on the projection (see Fig. 19a).

Using the inspector, the user sees that t-SNE is rather good at identifying regions in the data set. Some pieces of manifolds (e.g., of the digit 1 at the center of the projection, or 7 in Cluster 12, see Fig. 19b) are picked up as well. However, maybe due to the high variability of digit drawings, there is no clear boundaries between clusters. The visualization could thus be improved. In addition, the HD clustering performs badly, with initially 13% as Rand index.

Some clusters are clearly irrelevant, because their elements are scattered almost evenly in the visualization. Inspection confirms that their elements are rather unrelated (e.g., Clusters 1, 2, 8, 9 and 15, see Fig. 20b).

The user first focuses on the center of the projection, that looks more densely populated. A scope is thus defined by the selection of Clusters 11 and 3. Interestingly, despite a bad general performance of the clustering, these two specific clusters have captured a pattern absent from the projection: they seem to delimit a region for the digit 1, that was not clear at all from the sole projection. The user wants to emphasize this pattern. First, Clusters 11 and 3 are trivially merged via a simple legend edit. He or she then

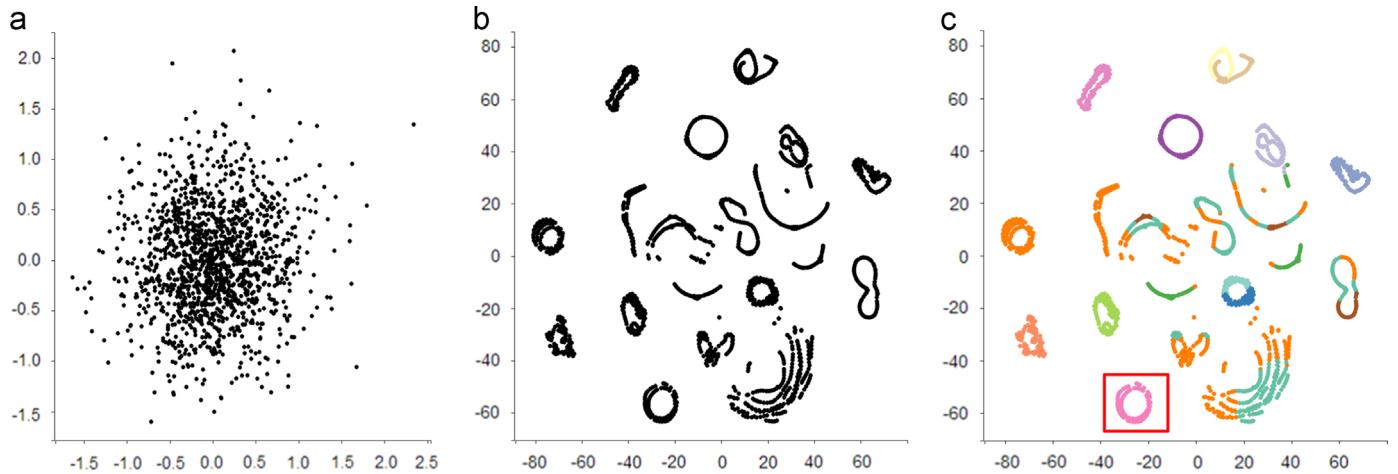


Fig. 14. (a) Initial configuration of t-SNE. (b) 2D projection after 1000 t-SNE iterations. (c) Overlay of the spectral clustering result. Cluster 14 is highlighted in red. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

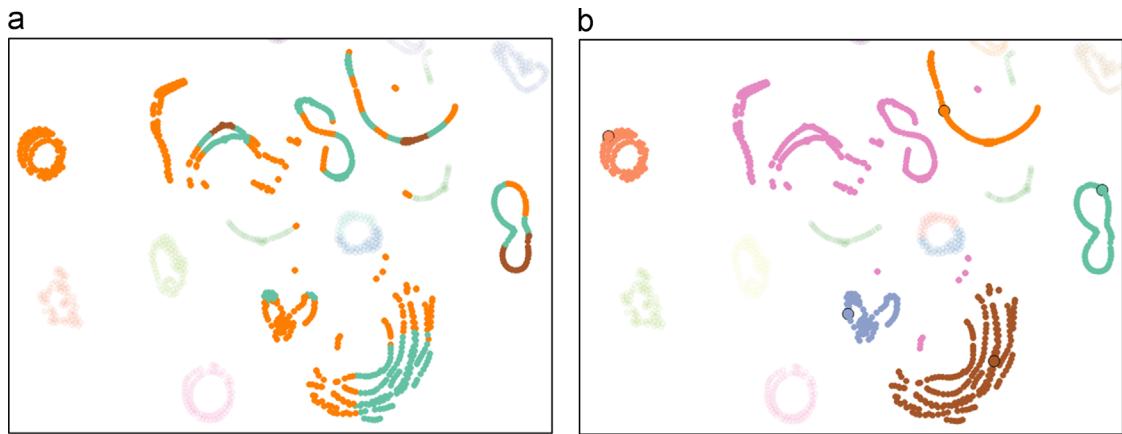


Fig. 15. (a) Highlight on Clusters 12 (orange), 13 (brown) and 15 (green). (b) Resulting clusters after isolation of manifolds from the highlighted seeds. The remaining elements in the scope become Cluster 18 (in purple). (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

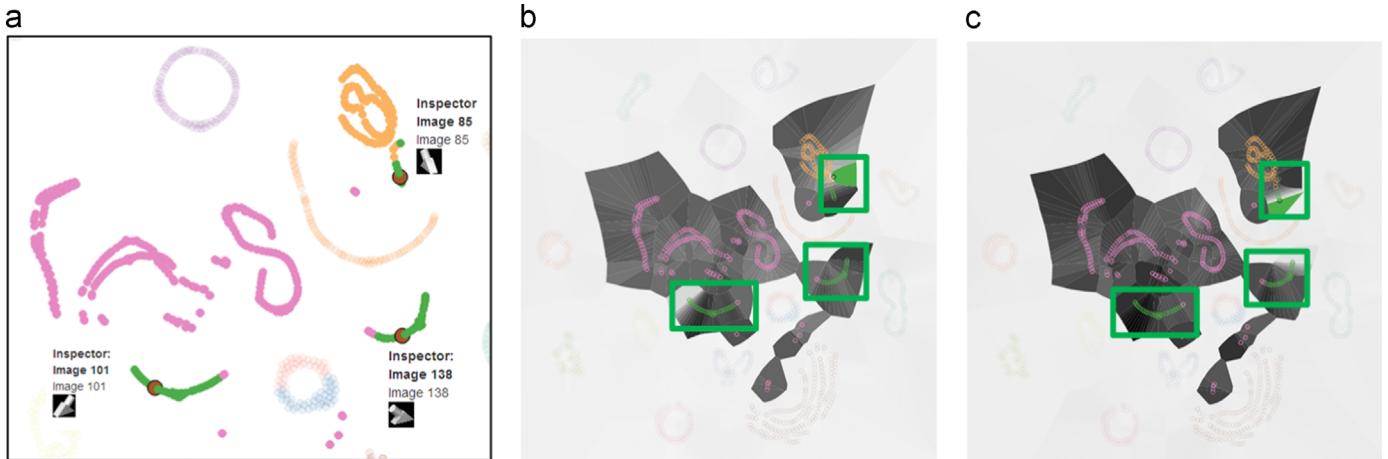


Fig. 16. (a) Scope containing Clusters 9 (orange), 10 (green) and 18 (purple). For each of the three components of Cluster 10, we attach the respective data inspector to one of its elements, highlighted in brown. (b) and (c) ProxiViz view of the scope, for two reference (hovered) elements, with Cluster 10 components highlighted in green. The scope extent is indicated using alpha blending. The Voronoi cells of non-hovered elements in the scope are mapped with gray shades, indicating the distance in the HD space w.r.t. the reference element. A tear occurs if an element is far to the reference element in the projection, but close in the HD space as indicated by the gray shades. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

applies a spreading transform, with a sharpness of 5, to this new cluster against all others (Fig. 20a).

Judging that most of the clustering appears as visual clutter, the user then focuses on Clusters 1, 2, 8, 9 and 15. Those exhibit

a poor locality w.r.t. the projection. The user defines a scope by selecting these clusters. He or she uses the inspector, and the marginal distribution of the scope in the visualization, to define the seeds for four rough regions. Next, the seeds are diffused

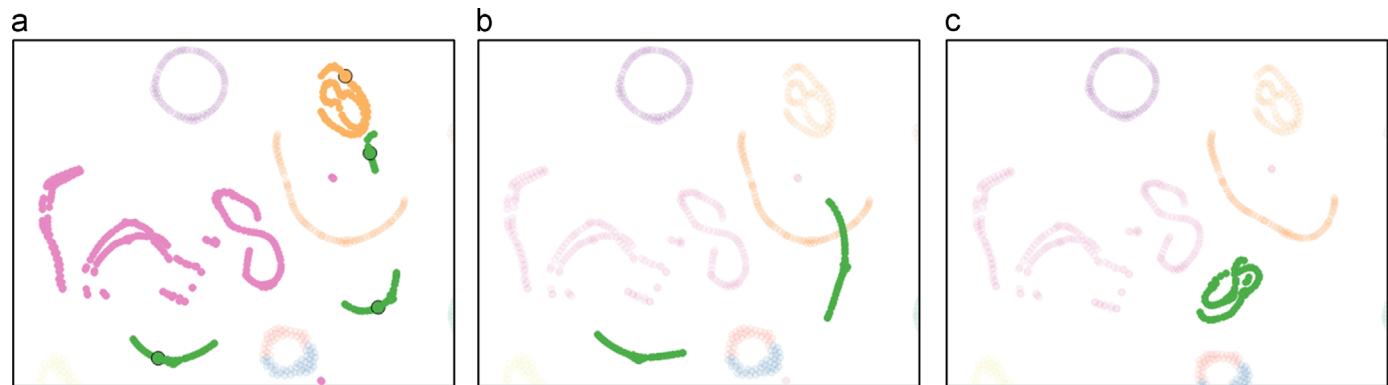


Fig. 17. (a) Cluster 10 after correction of the clustering errors using the label diffusion tools from three seeds. (b) Result after 50 scope-restricted t-SNE steps, without a dissimilarity transform. (c) Resulting projection after a gathering transform and 50 more scope-restricted steps.

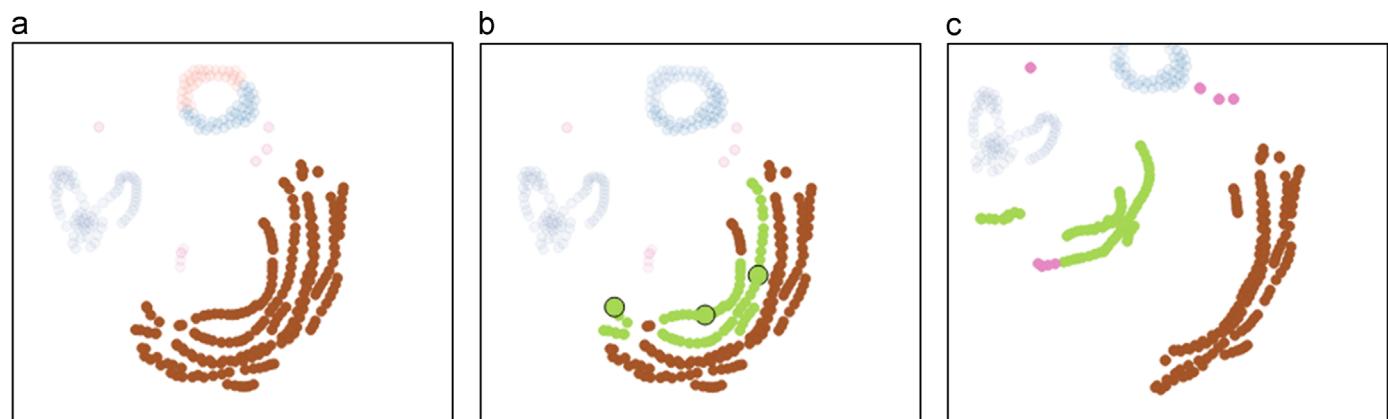


Fig. 18. (a) Highlight on Cluster 13. (b) Isolation of a subset of the cluster using three seeds. (c) Resulting projection and clustering, with the clearly separated Clusters 13 and 19. An outlying subset of Cluster 18 was spontaneously gathered to the new cluster.

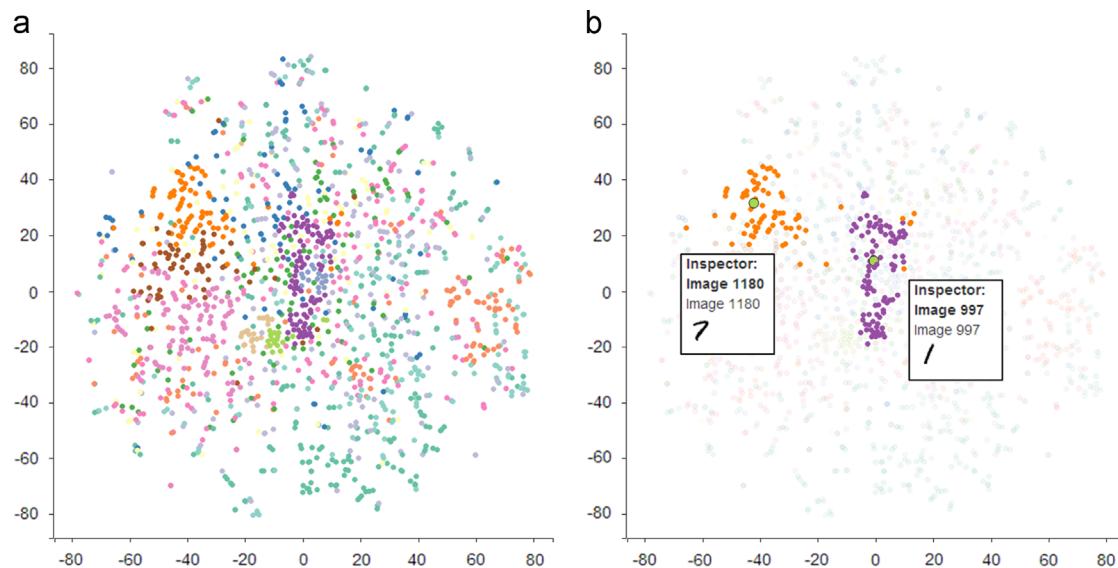


Fig. 19. (a) t-SNE projection of the MNIST image collection, with the spectral clustering result overlaid as categorical colors. (b) Cluster 11 (purple, in the center) contains almost exclusively the digit 1. Cluster 12 (orange, on the left) is mostly populated by samples of the digit 7. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

using the propagation tool, leading to more localized clusters (see Fig. 20c).

The user then notices Clusters 12 and 13, with both a more densely populated part and few outliers. Beyond their apparent proximity, inspection reveals they are both related to the same

ground truth digit (7). After trivially merging them with a legend edit, the user isolates the central manifold from the outliers with a single seed selection. The new merged cluster takes the label 12, and the remainder populates its own small *residual* cluster 13 (Fig. 21a and b).

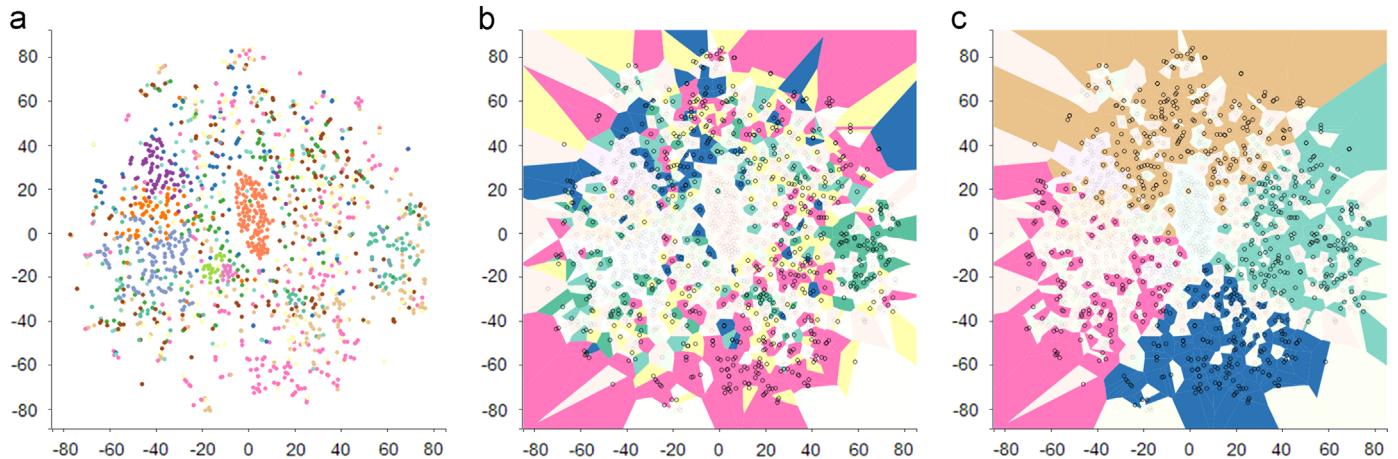


Fig. 20. (a) 2D projection after the trivial merger of Clusters 11 and 3. The spreading transform emphasizes the boundaries of the new cluster (center of the projection). (b) Highlight of Clusters 1, 2, 8, 9 and 15. The Voronoi visualization emphasizes the distribution in the 2D projection. (c) The same scope, after propagation of 4 seeds.

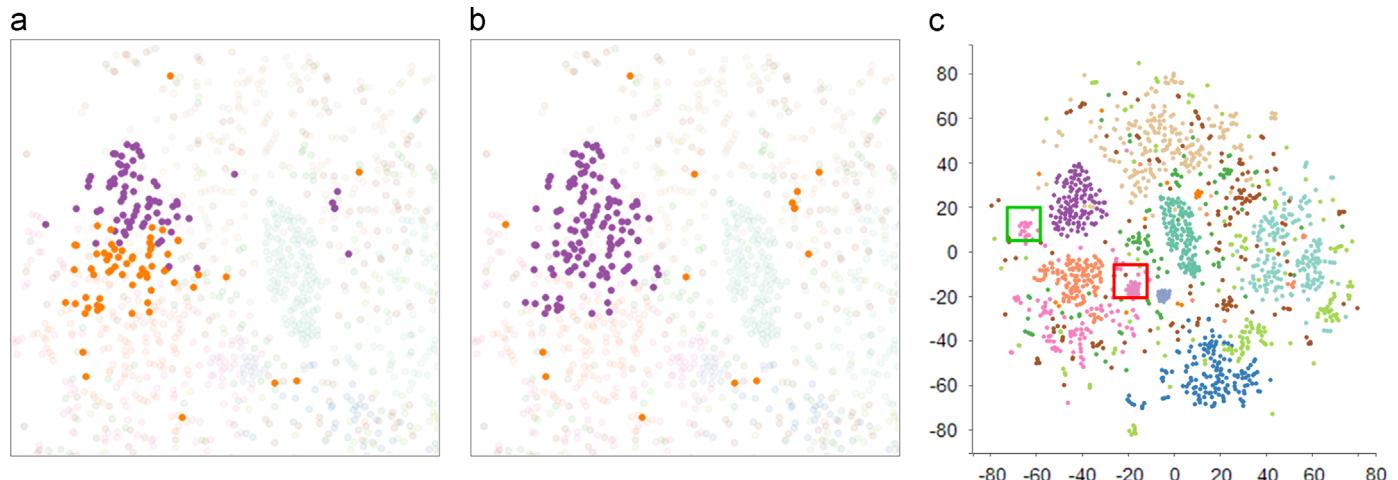


Fig. 21. (a) Clusters 12 (purple) and 13 (orange) before isolating the central manifold. (b) The central manifold becomes Cluster 12 (purple), and the remainder populates Cluster 13 (orange). (c) Visualization after applying a general spreading transform, and 50 t-SNE iterations. Cluster 6, highlighted in red, contains exclusively 1 digits, and the subset of Cluster 2 highlighted in green gathers samples of the digit 4. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

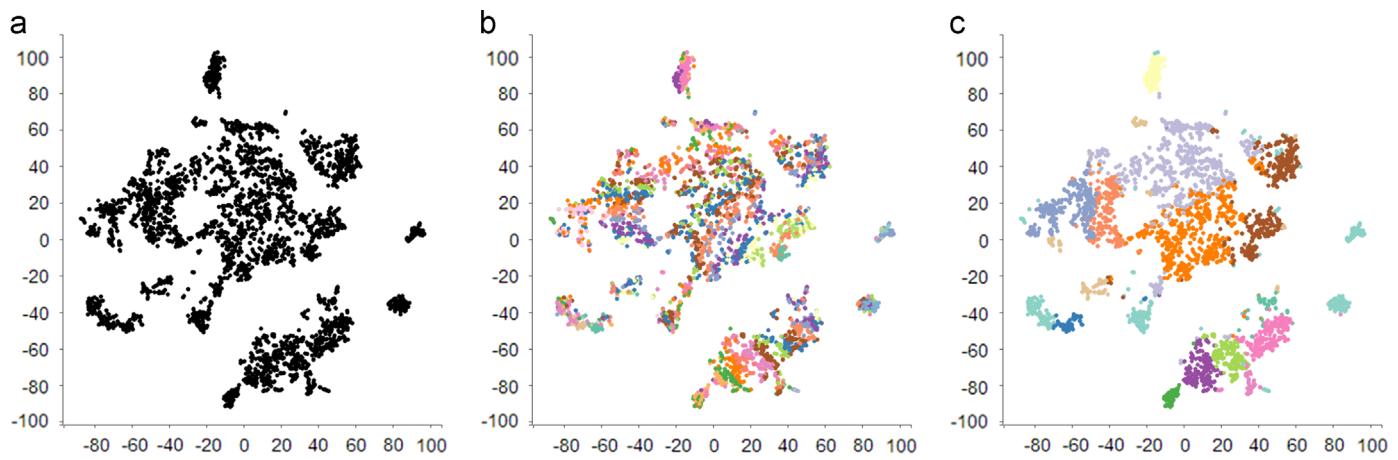


Fig. 22. (a) t-SNE projection of the *biogas* data set. (b) Overlay of the *k*-means clustering results. (c) Overlay of the spectral clustering results.

Finally, to emphasize the borders between the current clusters, the user applies a general spread transform (i.e., of all clusters w.r.t. all others) with a significantly high sharpness (10). After 50 iterations, the visualization looks much clearer, with a stronger emphasis given to manifolds and potential outliers (Fig. 21c).

Actually, this operation emphasized patterns unnoticed beforehand, e.g., Cluster 6, and dense regions in Cluster 2 (see Fig. 21c). The few interactions led to a significant improvement of the Rand index, now reaching 20%. This leaves an important margin for progression though.

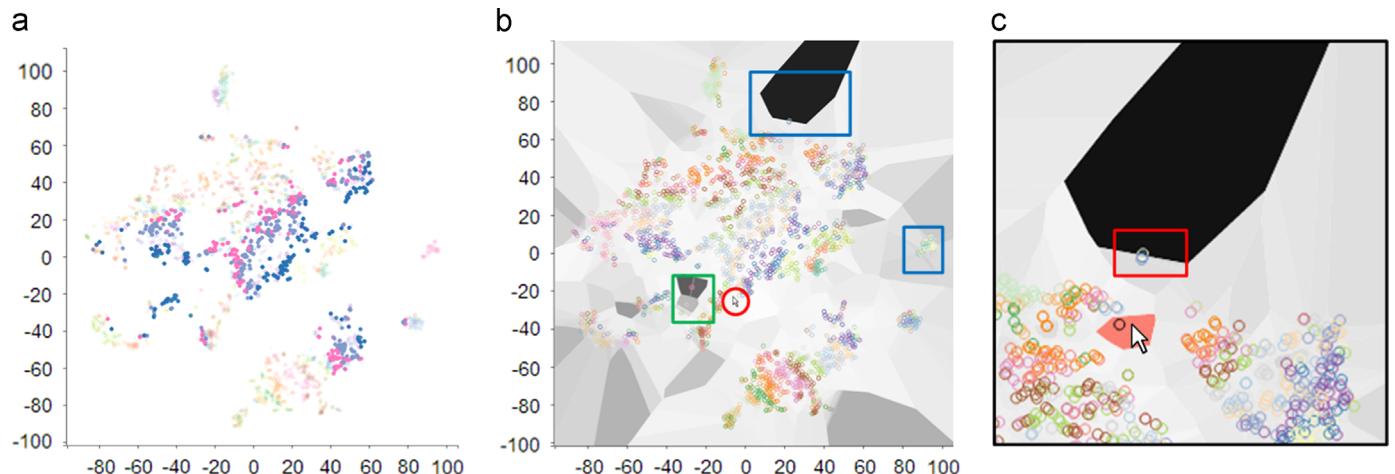


Fig. 23. (a) Highlight of Clusters 1 (blue), 4 (light blue) and 7 (pink). (b) ProxiViz distance pattern when hovering over the element highlighted by a red circle. Cluster 14 is highlighted in green. Cluster 8, highlighted in blue, contains the component on the right, and the outlying element on top of the visualization. (c) Among the four closely mapped elements (highlighted in red), from his dark shade, only one is clearly standing out. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

8.3. Biogas

Similar to the MNIST scenario, the user first loads a cached t-SNE projection of the *biogas* data (Fig. 22a). For this data set, two clustering results are available: a *k*-means output of 30 clusters, handed by the experts, and a spectral clustering output with the parametrization used in the two previous sections (15 clusters). Both results are overlaid on the 2D projection as categorical colors, respectively in Fig. 22b and c.

Using the data loading controls to compare the two clustering results, the user immediately notices a clear distinction between them. The spectral clustering results are quite closely related to the manifolds found by t-SNE. *k*-means on the other hand, finds a different kind of patterns, with clusters splitting across components (see Fig. 23a). As the experts provided the *k*-means results, the user chooses to use them as a starting point. The user expects to use his or her findings as a support for further interactions with the experts.

Before actually amending the visualization, the user wants to inspect the projection using the ProxiViz tool along with the Voronoi diagram visualization (see Fig. 23b). The user immediately identifies some outlying patterns. Some of the elements have a significantly darker shade irrespectively of the hovered element. This means that these elements are far from every other element in the HD space. Quoting the experts, such outliers are likely to be representing viral DNA. With its three elements clearly standing out, Cluster 14 (see Fig. 23b) has already been identified as such using parallel coordinates. The user found another pattern, more subtle, that could serve expert hypothesis formulation. With its dark shade, an element looks very clearly as an outlier. However, it was mapped very closely to three other elements that clearly do not look as such (Fig. 23c). The user also notes that the suspect element belongs to Cluster 8, that looks a bit outcast in the visualization, as confirmed by the ProxiViz view (Fig. 23b). Cluster 8 and the three elements mapped close to the suspect element are interesting cues for further analysis. The user is tempted to gather the components of Cluster 8. However, Fig. 23c highlights the heterogeneity of its elements. The chances thus are high that their position is due to a subtle equilibrium with many other elements in the 2D projection. In this context, gathering components in a restricted scope would require a very high sharpness, and could be widely disruptive.

The user next concentrates on Cluster 1. He or she wants to use the label diffusion and dissimilarity transform tools to regroup its components (see Fig. 24a). After having defined Cluster 1 as the current scope, the user adjusts the isolation cut control. This reflects the spread he or

she wants to reduce (Fig. 24b). No seed has been selected before, so the isolation effectively defines all components under the cut as new, temporary clusters (see Section 5.2). The user then gathers this new scope with a moderate sharpness (5). The temporary clusters that served as a cue for the gathering transform are then trivially restored as a single cluster (see Fig. 25a). The user is not satisfied with the resulting visualization though. He or she would like to have the new Cluster 1 more clearly separated from its neighbors, Clusters 4 and 11. Using spread transforms parametrized with a medium sharpness (10), the visualization shown in Fig. 25b is eventually obtained.

This scenario is part of a preliminary stage to the exploratory analysis. The user notes that the *biogas* data set lacks a summary representation for its elements. Plugging such representations in the data inspection tool would give a crucial cue for further insight. The design of such a representation is thus to be shortly discussed with the experts.

9. Discussion

In the context of interactive systems, the computational complexity of the involved tools is critical. Four potential bottlenecks exist in Cluster Sculptor:

- the computation of the Minimum Spanning Tree (MST),
- the interactive update of the MST cut,
- the dissimilarity transform,
- the t-SNE updates.

The graph underlying t-SNE is complete, causing MST computation to be $O(N^2)$. This is not excessively harmful, as the MST has to be updated only after a dissimilarity transform. Also, using an optimal MST is not critical to the application. Therefore, increased speed could be obtained by filtering the highest range of values from the graph implicit to \mathbf{D} .

Updating the cut is linear w.r.t. the number of edges in the MST, thus linear w.r.t. the number of elements in the scope, as the MST of a graph defined over N elements has $N - 1$ edges [37,38]. This operation is generally quite fast, but can become slow if no scope is defined, and the number of elements exceeds 5000. However, the current implementation is naive, and dramatic improvements could be made using a properly sorted data structure.

The dissimilarity transform is $O(N^2)$ in the worst case. However, this step just applies P_{β} on graph edges in the selected scope

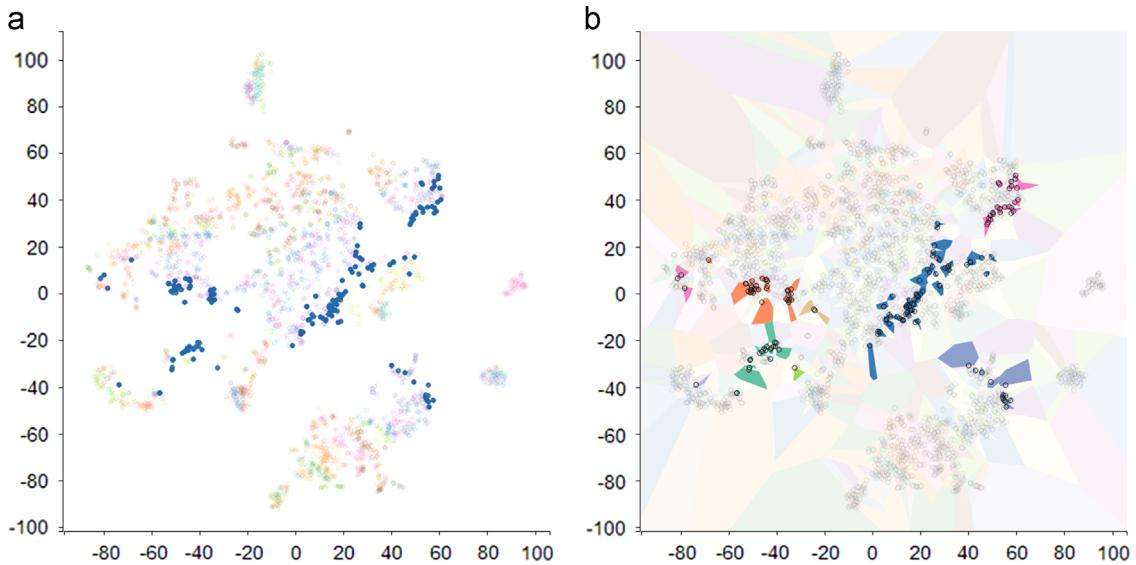


Fig. 24. (a) Scope containing Cluster 1 (blue). (b) Voronoi visualization of the components isolated by the MST cut. (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

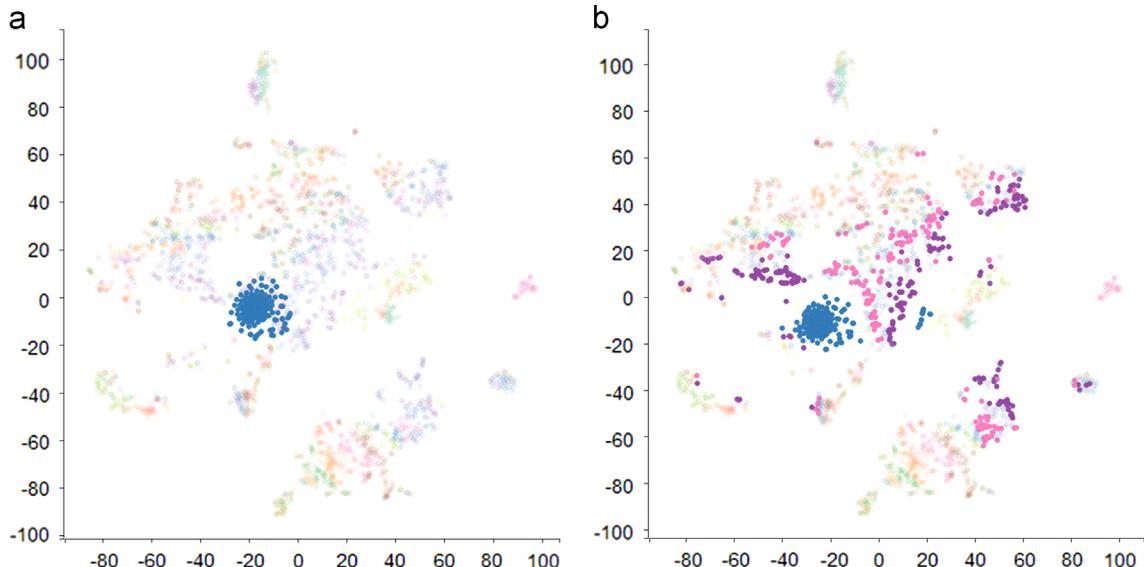


Fig. 25. (a) Cluster 1 (blue) after the gathering transform, and 200 t-SNE iterations. (b) Cluster 1 after being spread from conflicting neighbors (Clusters 4 (pink) and 11 (purple)). (For interpretation of the references to color in this figure caption, the reader is referred to the web version of this paper.)

(see Section 6). Consequently, each operation is considerably faster than those of the MST computation. For a higher interactivity, in case of very large data sets, piecewise computation could also be considered.

t-SNE updates are $O(N^2)$ [14]. Unlike dissimilarity transforms, they occur frequently, and their speed is the basis for the interactivity of the system. A user will perceive the succession of updates as a smooth move only if the frequency of updates is sufficiently high. Some studies argue that a system would be perceived as interactive only if the response time is under 100 ms [12]. Unfortunately, in the current implementation, updating the position of 2000 elements in the 2D projection takes approximately 1 s. However, we did not investigate improvements to the baseline t-SNE steps, such as a random-walk approximation [14], or more recently $O(N \log N)$ variants based on the Barnes-Hut algorithm [41,42]. Using such improvements, we could for example easily process the complete biogas data set (see Section 8).

A major issue in our system is to keep the user engaged when making changes to the projection. Heer and Robertson found that

careful animation design has significant advantages for graphical perception of changes in the context of statistical plots [43]. When applying a dissimilarity transform, the visual tracking of objects is facilitated by the physical metaphor t-SNE follows (see Section 3). However, if the frame rate of t-SNE updates is too low, the sequence may not be perceived as a smooth move. This issue could be potentially alleviated by using interpolation and animation features available in d3.js [44]. We already use this library widely in the Cluster Sculptor prototype.

In this paper, we emphasized the use of the label diffusion tools, instead of classical lasso or rectangular selectors. In the current prototype state, this is the only supported selection mode. This design choice was motivated by our intuition, and should eventually be validated by a proper user study. In the meantime, classical selection modes could also be implemented. We sense that the most efficient tool should depend on the properties of the selection to be performed.

We purposely chose to apply dissimilarity updates on the basis of the set of labels in the scope. This is unlike specific pairwise

constraints, such as seen in metric learning approaches [45,30]. We found that these methods require too high amounts of user-specified constraints to be really effective, and wanted to spare user efforts. For example, Martin et al. require approximately 30 pairwise constraints for their technique to be effective [45]. The approach of Basu et al. requires at least 20% of the collection to be labeled [31]. In Dis-Function [30], the authors choose to bias their objective function and give higher importance to the user specified constraints. We alternatively aimed at implementing smooth, non-disruptive updates to the visualization, in a more robust fashion. Yet, a comparative analysis, or a way of combining the approaches should be investigated.

Our approach has a strong weakness: unlike the metric learning approaches, it is not currently able to generalize to new data points. This can be problematic in the context of data streams and further work is needed to support them.

Another important point regards the updated dissimilarity matrix, \mathbf{D} . It is explicitly stated as a dissimilarity matrix as its transform by P_{β} may invalidate its initial Euclidean metricity. Fortunately, t-SNE is not limited by the metricity of these input values, and is able to cope with any valid dissimilarity matrix.

In principle our approach is DR-technique agnostic and could be thus combined with any DR technique, as long as the above mentioned concerns are addressed. The choice of technique also depends on various factors. The format of the input (e.g., distance matrix or coordinates in original space) is one of these. The trade-off between computational complexity and projection quality can also be influential [46–48].

Finally, related work such as iVisClustering [4] and Dis-Function [30] demonstrate the use of multiple coordinated views. This can help in gaining a better insight and understanding of the data and model at hand. This direction has not been explored in our work, though some needs have been identified in Section 8.3 (e.g., an adapted inspector view or complementary insight by parallel coordinates). The system could also benefit from the dual use of both views introduced in the paper (i.e., scatterplot view and Voronoi diagram view). Linking these views would raise potential for new and interesting features, such as the inspection of distance patterns in the HD space with joint dynamical emphasis of clusters. Complementary interactions, that go beyond mere juxtaposition of views, should be investigated.

10. Conclusion and perspectives

This paper described Cluster Sculptor, a novel interactive clustering system. Our system takes clustering results as an input, and extensively uses the t-SNE projection technique. Its iterative nature, and the physical metaphor underlying it are a key to our approach. Its contributed components pertain essentially to a framework of interactive tools, i.e., label diffusion and dissimilarity transform. Cluster Sculptor also grounds in the ability of t-SNE to adapt smoothly to discontinuous dissimilarity updates. As such, we are able to prevent major distortions in the user's mental map. Detailed experimental scenarios, using real data sets, showed how a user could combine the diverse features of the system to amend clustering results, and the associated 2D projection.

This work opens to numerous perspectives. The tools described in the paper are prototypical, and could be improved in many aspects. For example, we included the ProxiViz tool [49] to answer the most urgent needs. However, we did not account for ProxiLens, the most recent evolution of the tool [50]. It features an interactive lens, that animates tears and false neighborhoods when hovering over any point in the visualization. Its integration could be considered in future work.

We are aware that the system lacks basic HCI features, such as a history of interactions, or at least undo and redo features. Their

implementation would help reduce user frustration, render the system more usable, and to some extent give some further insight of the data via the post-hoc navigation in interaction sequences. It is also clear that in its current, prototypical, state, the tool is not usable by someone with no knowledge in machine learning and DR methods. We showed the feasibility of the approach, and left the simplification of the interface for future work.

The discussion has shown openings for a better scalability. If very large and dense data sets are targeted, Cluster Sculptor could be supplemented by tools for interacting with dense patterns, such as a zooming feature, or a fish-eye plugin [51].

In the experimental section, numerous examples of MST cuts are given. As a complement to the isolation cut parameterization, visualizing and interacting directly with these graphs could give increased control to the user.

The dissimilarity transform is currently applied to all clusters in the current scope. For example, when gathering the components underlying a cluster, or spreading a cluster w.r.t. a subset of its neighbors, this design leads to tedious relabeling. A set of predefined sequences could be derived, e.g. a gathering operator including propagation and isolation steps.

The generalization of Cluster Sculptor to more data types (e.g., categorical) and DR techniques could be investigated in future work. Also, the relationship, and potentially complementarity, with methods based on learning a distance function using pairwise constraints, is an interesting direction of research. Detecting the local relevance of HD features would be an interesting perspective to limit the disruptive tendency of classical metric learning approaches. Accounting for the latter aspects could widen the applicability of the system, e.g. to heterogeneous data types, or lessening the importance of noisy features.

Finally, the experimental scenarios involve a variety of low-level tasks (e.g., isolating a cluster with a MST cut). As evoked in the experimental section, we chose to put an emphasis on scenarios of use: but these tasks should certainly be evaluated on a unitary basis in the context of a classical user study.

Acknowledgments

We would like to warmly thank the anonymous referees for their extensive reviews. They highlighted important points that largely contributed to the improvement of this work.

References

- [1] M.B. Eisen, P.T. Spellman, P.O. Brown, D. Botstein, *Cluster analysis and display of genome-wide expression patterns*, Proc. Natl. Acad. Sci. 95 (25) (1998) 14863–14868.
- [2] T. Schreck, J. Bernard, T. von Landesberger, J. Kohlhammer, Visual cluster analysis of trajectory data with interactive kohonen maps, Inf. Vis. 8 (1) (2009) 14–29. <http://dx.doi.org/10.1057/ivs.2008.29> arXiv:<http://ivi.sagepub.com/content/8/1/14.full.pdf+html> URL: <http://ivi.sagepub.com/content/8/1/14.abstract>.
- [3] G. Andrienko, N. Andrienko, S. Rinzivillo, M. Nanni, D. Pedreschi, F. Gianotti, Interactive visual clustering of large collections of trajectories, in: IEEE Symposium on Visual Analytics Science and Technology, 2009, pp. 3–10.
- [4] H. Lee, J. Kihm, J. Choo, J. Stasko, H. Park, iVisClustering: an interactive visual document clustering via topic modeling, in: Computer Graphics Forum, vol. 31, Wiley Online Library, 2012, pp. 1155–1164.
- [5] A.K. Jain, *Data clustering: 50 years beyond k-means*, Pattern Recognit. Lett. 31 (8) (2010) 651–666.
- [6] A.Y. Ng, M.I. Jordan, Y. Weiss, On Spectral Clustering: Analysis and an Algorithm, NIPS.
- [7] C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, New York, 2006.
- [8] S. Léspinats, M. Aupetit, Classimap: a supervised mapping technique for decision support, in: Proceedings of the EuroVis Workshop on Visual Analytics using Multidimensional Projections, 2013, pp. 17–20. <http://dx.doi.org/10.2312/PE.VAMP.VAMP2013.017-020> URL: <http://diglib.eg.org/EG/DL/PE/VAMP/VAMP2013/017-020.pdf>.
- [9] A. Karatzoglou, A. Smola, K. Hornik, kernlab (R package) (2013).

- [10] I. Fodor, A Survey of Dimension Reduction Techniques, Technical Report, Lawrence Livermore National Laboratory (2002).
- [11] L. Van der Maaten, E. Postma, H. Van Den Herik, Dimensionality Reduction: A Comparative Review, Technical Report, Tilburg University (2009).
- [12] C. Ware, *Information Visualization: Perception for Design*, Elsevier, Morgan Kaufmann, San Francisco, 2004.
- [13] P. Bruneau, F. Picarougne, M. Gelgon, Interactive unsupervised classification and visualization for browsing an image collection, *Pattern Recognit.* (2010) 485–493.
- [14] L. van der Maaten, G. Hinton, Visualizing data using t-SNE, *J. Mach. Learn. Res.* 9 (2008) 2579–2605.
- [15] S.A. Nene, S.K. Nayar, H. Murase, Columbia object image library (coil-20), Technical Report CUUCS-005-96 (1996).
- [16] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, *Proc. IEEE* 86 (11) (1998) 2278–2324.
- [17] J. Seo, B. Shneiderman, Interactively exploring hierarchical clustering results [gene identification], *Computer* 35 (7) (2002) 80–86. <http://dx.doi.org/10.1109/MC.2002.1016905>, URL: <<http://dx.doi.org/10.1109/MC.2002.1016905>>.
- [18] C. Turkay, J. Parulek, N. Reuter, H. Hauser, Integrating cluster formation and cluster evaluation in interactive visual analysis, in: Proceedings of Spring Conference on Computer Graphics.
- [19] S. Rinzivillo, D. Pedreschi, M. Nanni, F. Giannotti, N. Andrienko, G. Andrienko, Visually driven analysis of movement data by progressive clustering, *Inf. Vis.* 7 (3–4) (2008) 225–239. <http://dx.doi.org/10.1057/palgrave.ivs.9500183>, URL: <<http://dx.doi.org/10.1057/palgrave.ivs.9500183>>.
- [20] A. Kapoor, B. Lee, D. Tan, E. Horvitz, Interactive optimization for steering machine classification, in: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '10, ACM, New York, NY, USA, 2010, pp. 1343–1352. <http://dx.doi.org/10.1145/1753326.1753529>, URL: <<http://doi.acm.org/10.1145/1753326.1753529>>.
- [21] B. Broeksema, A.C. Telea, T. Baudel, Visual analysis of multi-dimensional categorical data sets, *Comput. Graph. Forum* 32 (8) (2013) 158–169. <http://dx.doi.org/10.1111/cgf.12194>, URL: <<http://dx.doi.org/10.1111/cgf.12194>>.
- [22] J. Philippeau, J. Pinquier, P. Joly, J. Carrive, Dynamic organization of audiovisual database using a user-defined similarity measure based on low-level features, in: IEEE International Conference on Image Processing, 2008, pp. 33–36.
- [23] G.M.H. Mamani, F.M. Fatore, L.G. Nonato, F.V. Paulovich, User-driven feature space transformation, *Comput. Graph. Forum* 32 (3) (2013) 291–299.
- [24] M. Aupetit, Visualizing distortions and recovering topology in continuous projection techniques, *Neurocomputing* (2007) 1304–1330.
- [25] I.T. Jolliffe, *Principal Component Analysis*, Springer, New York, 1986.
- [26] P. Desmartins, J. Héault, Curvilinear component analysis: a self-organising neural network for non-linear mapping of data sets, *IEEE Trans. Neural Netw.* (1997) 148–154.
- [27] N. Heulot, M. Aupetit, J.-D. Fekete, Proxiziv: an interactive visualization technique to overcome multidimensional scaling artifacts, in: Proc. IEEE InfoVis, poster, 2012.
- [28] D. Shepard, A two-dimensional interpolation function for irregularly-spaced data, in: Proceedings of the Twenty-third ACM National Conference, 1968, pp. 517–524.
- [29] L. Martin, M. Exbrayat, G. Cleuziou, F. Moal, Interactive and progressive constraint definition for dimensionality reduction and visualization, in: G.R.F. Guillet, D. Zighed (Eds.), *Advances in Knowledge Discovery and Management (AKDM-2)*, Studies in Computational Intelligence, vol. 2, Springer, 2012, pp. 121–136, URL: <<http://hal.archives-ouvertes.fr/hal-00595035>>.
- [30] E.T. Brown, J. Liu, C.E. Brodley, R. Chang, Dis-function: learning distance functions interactively, in: IEEE Conference on Visual Analytics Science and Technology, 2012, pp. 83–92.
- [31] S. Basu, A. Banerjee, R. Mooney, Semi-supervised clustering by seeding, in: Proceedings of Nineteenth International Conference on Machine Learning.
- [32] X. Zhu, Z. Ghahramani, Learning from Labeled and Unlabeled Data with Label Propagation, Technical Report CMU-CALD-02-107, Carnegie Mellon University (2002).
- [33] P. Bruneau, B. Otjacques, A proposition of interactive visual clustering system, EuroVis 2013 Workshop on Visual Analytics using Multidimensional Projections.
- [34] A. Gisbrecht, B. Hammer, B. Mokbel, A. Sczryba, Nonlinear dimensionality reduction for cluster identification in metagenomic samples, in: Seventeenth International Conference on Information Visualisation (IV), 2013, pp. 174–179.
- [35] S. Tilkov, S. Vinoski, Node.js: using javascript to build high-performance network programs, *IEEE Internet Comput.* 14 (6) (2010) 80–83.
- [36] M. Harrower, C.A. Brewer, ColorBrewer.org: an online tool for selecting colour schemes for maps, *Cartograph. J.* 40 (1) (2003) 27–37.
- [37] J.B. Kruskal, On the shortest spanning subtree of a graph and the traveling salesman problem, in: Proceedings of the American Mathematical Society, vol. 7, 1956, pp. 48–50.
- [38] R.C. Prim, Shortest connection networks and some generalizations, *Bell Syst. Tech. J.* (1957) 1389–1401.
- [39] J. Donaldson, T-distributed Stochastic Neighbor Embedding for R (R package) (2012). URL: <<http://cran.r-project.org/web/packages/tSNE/index.html>>.
- [40] L. Hubert, P. Arabie, Comparing partitions, *J. Classif.* 2 (1985) 193–218.
- [41] L. Van der Maaten, Barnes-hut-sne, Technical report, Delft University of Technology (2013).
- [42] Z. Yang, J. Peltonen, S. Kaski, Scalable optimization of neighbor embedding for visualization, in: Proceedings of the Thirtieth International Conference on Machine Learning, 2013, pp. 127–135.
- [43] J. Heer, G.G. Robertson, Animated transitions in statistical data graphics, *IEEE Trans. Vis. Comput. Graph.* 13 (6) (2007) 1240–1247.
- [44] M. Bostock, Data driven documents, URL: <<http://d3js.org>> (2013).
- [45] L. Martin, M. Exbrayat, G. Cleuziou, F. Moal, Interactive and progressive constraint definition for dimensionality reduction and visualization, *Adv. Knowl. Discov. Manag. Stud. Comput. Intell.* 398 (2012) 121–136.
- [46] I. Borg, *Modern Multidimensional Scaling: Theory and Applications*, Springer, New York, 2005.
- [47] F.V. Paulovich, C.T. Silva, L.G. Nonato, Two-phase mapping for projecting massive data sets, *IEEE Trans. Vis. Comput. Graph.* 16 (6) (2010) 1281–1290.
- [48] S.L. France, J. Carroll, Two-way multidimensional scaling: a review, *IEEE Trans. Syst. Man Cybern. C Appl. Rev.* 41 (5) (2011) 644–661.
- [49] N. Heulot, M. Aupetit, J.-D. Fekete, Évaluation de proxiviz pour la fouille visuelle de données multidimensionnelles, Atelier Fouille Visuelle de Données: méthodologie et évaluation, in: ECG 2012.
- [50] N. Heulot, M. Aupetit, J.-D. Fekete, Proxilens: interactive exploration of high-dimensional data using projections, in: Proceedings of the EuroVis Workshop on Visual Analytics using Multidimensional Projections, 2013, pp. 11–15.
- [51] M. Sarkar, M.H. Brown, Graphical fisheye views of graphs, in: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, 1992, pp. 83–91.



Pierrick Bruneau holds a M.Sc degree in computer sciences, received from PolytechNantes (France) in 2007. He also obtained a Ph.D. from the University of Nantes in 2010, and conducted postdoctoral research at CEA LIST (Saclay, France). He is now a researcher at the Informatics, Systems and Collaboration department of the Gabriel Lippmann Public Research Center (Belvaux, Luxembourg). His research is focused on data mining, applied statistics, and visual analytics.



Philippe Pinheiro holds a M.Sc degree in computer sciences, received from ISIAL Nancy (France) in 2001. He is a senior Java developer and worked three years for an international insurance company. He is now a researcher at the Informatics, Systems and Collaboration department of the Gabriel Lippmann Public Research Center (Belvaux, Luxembourg). His work is focused on applied research in visual analytics and data mining.



Bertjan Broeksema holds a M.Sc degree in computer science, received from the university of Groningen in 2010. He expects to obtain his Ph.D from the university of Groningen early 2014, as a result of his research at IBM France on visual analytics solutions for decision management systems. He is now a researcher at the Informatics, Systems and Collaboration department of the Gabriel Lippmann Public Research Center (Belvaux, Luxembourg). His research is focused on visual analytics, information visualization and data mining.



Benoît Otjacques received his "Ingénieur civil" degree in 1992 from the University of Louvain, Belgium and his Ph.D in Computer Science from the University of Namur, Belgium in 2009. He also received a Master in Law and Management of ICT (University of Namur, 1993), a Master in Innovation Management (University of Mons, Belgium, 2000), and a Research Master in Modeling and Simulation in Architecture (University of Nancy, France, 2004). Since 2002, he is deputy scientific director of the Informatics, Systems and Collaboration (ISC) Department of the Gabriel Lippmann Public Research Center, Luxembourg. His main research topics are information visualization, visual analytics and human-computer interaction.