

Visualizing Distance Preservation

As outlined in Chapter 2, multidimensional projections are effective tools for exploring large sets of high-dimensional observations, especially when the questions of interest regard the detection of groups of similar observations or isolated outlier observations. However, as the same chapter has explained, projection techniques are challenged by the inability of faithfully preserving the so-called structure of the high-dimensional data, such as inter-point distances or point neighbors. If such aspects of the data structure are not well preserved – or, alternatively, if the lack of preservation is not clearly shown to the user, the visualization may convey wrong insights in the data.

In this chapter, we address the first part of the above challenge, namely measuring and displaying *distance preservation errors* in multidimensional projections. In detail, we present a set of metrics that aim to address the following questions for 2D projections:

- How is the projection error spread over the 2D space?
- How to find points that are close in 2D but far in n D?
- How to find points that are close in n D but far in 2D?
- How do the choice of projection algorithm and its parameter settings affect the above quality aspects?

To visualize the proposed metrics, we develop several space-filling techniques that visually scale to large datasets, offer a multiscale (or level-of-detail) view on the projection behavior, and do not require users to understand the internal formulation of dimensionality-reduction algorithm. We next use these visualizations to explore how five state-of-the-art dimensionality-reduction techniques behave, in terms of distance-preservation errors, when varying their parameters. Our endeavor of exploring projection errors is completed, next, by the work presented in Chapter 4, which addresses the related, but different, task of measuring and analysing neighborhood preservation errors for multidimensional projections.

3.1 Analysis Goals

Let us first recall the basic notations used to describe multidimensional projections. A multidimensional projection technique was modeled as a function $f : \mathbb{R}^n \times P \rightarrow \mathbb{R}^m$ that takes a dataset $D^n \subset \mathbb{R}^n$ and maps it to a lower dimensional dataset $D^m \subset \mathbb{R}^m$, $m < n$ (see Eqn. 2.1 and related text in Sec. 2.3). Here, P indicates the space of parameters of the projection technique f , which depends on the specific details of the algorithm underlying f .

A projection f should preserve the *structure* of the original space \mathbb{R}^n . This implies, as discussed in Chapter 2, a mix of distance and neighborhood preservations at various

scales and happens at different rates for different datasets, projection algorithms, and parameter values. One important task for which projections are used is to detect and reason about groups of close points in D^m . Unless users are sure that such groups of points are indeed also close in the original high-dimensional dataset D^n , the interpretation of the projection result can be misleading [8]. Hence, for such tasks, the perceived precision (or, in other words, quality) of a projection is intrinsically linked to its ability to preserve distances between points [159].

Our exploration goal can be thus refined as follows: Given any dimensionality-reduction (DR), or projection, algorithm (Eqn. 2.1), we aim to show how distance preservation is affected by choices of parameter values in P , highlighting aspects that can adversely affect the interpretation of the projected point set D^m . To simplify the discourse, we next consider $m = 2$, and that projections are drawn as scatterplots – the most common option for DR visualization. We identify the following aspects of interest in the exploration of distance preservation in multidimensional projections:

- A. False neighbors:** Take a point $\mathbf{p}_i \in D^n$ and its 2D projection $\mathbf{q}_i = f(\mathbf{p}_i)$. A necessary condition for distance preservation is that *all* points \mathbf{q}_j that are close to \mathbf{q}_i (in 2D) should be projections of points \mathbf{p}_j that are close to \mathbf{p}_i (in D^n). If not, *i.e.* we have a \mathbf{q}_j close to \mathbf{q}_i for which \mathbf{p}_j is not close to \mathbf{p}_i , the user wrongly infers from the projection that \mathbf{p}_j is close to \mathbf{p}_i . We call such a point j a *false neighbor*¹ of i .
- B. Missing neighbors:** The second necessary condition for distance preservation is that *all* \mathbf{p}_j that are close to \mathbf{p}_i (in D^n) project to points \mathbf{q}_j that are close to \mathbf{q}_i (in 2D). If not, *i.e.* we have a \mathbf{p}_j close to \mathbf{p}_i for which \mathbf{q}_j is not close to \mathbf{q}_i , the user will underestimate the set of points similar to point i . We call such a point j a *missing neighbor* of i .
- C. Groups:** A main goal of DR is to help users find groups of similar points, *e.g.* topics in a document set [95, 139] or classes of images in a database [55]. False and missing point neighbors generalize, for groups, to *false members* and *missing members* respectively. Given a group Γ of closely projected points, we aim to find if *all* points in Γ truly belong there (no false members), and if *all* points that belong to the topic described by Γ do indeed project in Γ (no missing members).
- D. Detail:** Aggregated local metrics such as [159, 8, 111, 79] can show, up to various extents, where missing or false neighbors occur. However, they do not directly show which are all such neighbors, for each projected point. Also, they do not explicitly address locating false and missing group members. We aim to provide interactive visual mechanisms to support these tasks on several levels of detail.

¹It is important to note that, although our discourse here uses terms such as false neighbors and missing neighbors, the projection-error visualization techniques discussed in this chapter focus on showing *distance*-preservation errors and not *neighborhood*-preservation errors, the latter being the subject of Chapter 4. We employ here terms that use the ‘neighbor’ concept simply because this term is more compact, and arguably more illustrative, than alternative distance-related terms.

We next propose several visualization methods to address the analysis goals outlined in this section.

3.2 Visualization Methods

As a running example, we use the Locally Affine Multidimensional Projection (LAMP) technique as projection method, with the default parameter settings given in [95], and as input the well-known 19-dimensional Segmentation dataset with 2300 points from [116, 95, 138, 141]. In this dataset, each point describes a randomly drawn 3x3 pixel-block from a set of 7 manually segmented outdoor images, by means of 19 statistical image attributes, such as color mean, standard deviation, and horizontal and vertical contrast. The aim of analysing this multidimensional dataset is to see how well extracted features group over observations, and thereby derive insights that can further on help in the construction of automatic classifiers for natural images. For completeness, we note that the task of classifier construction is outside our scope – our goal is strictly the exploration of the capability of projection techniques in displaying the similarity of high-dimensional observations, and errors involved in this display.

3.2.1 Preliminaries

To quantify the distance preservation issues in Sec. 3.1, we first define the projection error of point i vs a point $j \neq i$ as

$$e_{ij} = \frac{d^m(\mathbf{q}_i, \mathbf{q}_j)}{\max_{i,j} d^m(\mathbf{q}_i, \mathbf{q}_j)} - \frac{d^n(\mathbf{p}_i, \mathbf{p}_j)}{\max_{i,j} d^n(\mathbf{p}_i, \mathbf{p}_j)}. \quad (3.1)$$

We see that $e_{ij} \in [-1, 1]$. Negative errors indicate points whose projections are too close (thus, false neighbors). Positive errors indicate points whose projections are too far apart (thus, missing neighbors). Zero values indicate ‘good’ projections, which approximate optimally the distances in D^n .

Comparing Eqn. 3.1 with the aggregated normalized stress definition σ (Eqn. 2.2), we see similarities but also several differences. First, both errors are essentially differences between the distances in D^m and D^n between data points – the more these distances are different, the higher is the error. However, our error e_{ij} is not a strictly positive value as the (terms of the) stress function. This allows us to distinguish, on a more refined level, points that are projected too close from points that are projected too far away. Secondly, our error is normalized between -1 and 1 , which allows us to compare more easily different projections and/or different points in the same projection. In turn, both above aspects allow us to easily create several detail visualizations of projection errors. These are described next.

3.2.2 The Aggregated Error view

We first provide an overview of how the projection error spreads over an entire dataset, by computing for each point i the aggregate error

$$e_i^{agg} = \sum_{j \neq i} |e_{ij}|. \quad (3.2)$$

The value of e_i^{agg} gives the projection error of point i with respect to *all* other points. Low values of e^{agg} show points whose projections can be reliably compared with most other projections in terms of assessing similarity. These are good candidates for representatives in multilevel projection methods [55, 144, 29, 139]. Large values of e^{agg} show points that are badly placed with respect to most other points. These are good candidates for manual projection optimization [142, 138].

Fig. 3.1 (a) shows e^{agg} by color mapping its value on the 2D projected points, using a blue-yellow-red diverging colormap [74]. Brushing and zooming this image allows inspecting e^{agg} for individual points. However, given our goal of providing an overview first, we are actually not interested in *all* individual e^{agg} values, but rather to (a) find compact areas in the projection having similar e^{agg} values, (b) find outlier e^{agg} values in these areas (if any), and (c) see how e^{agg} globally varies across the projection. For this, we propose an image-based, space-filling visualization, as follows. Denote by $DT(\mathbf{x} \in \mathbb{R}^2) = \min_{\mathbf{q} \in D^m} \|\mathbf{q} - \mathbf{x}\|$ the so-called distance transform of the 2D point cloud D^m delivering, for any screen pixel \mathbf{x} , its distance to the closest point in D^m [39]. We then compute e^{agg} at every screen pixel \mathbf{x} as

$$e^{agg}(\mathbf{x}) = \frac{\sum_{\mathbf{q} \in N_\epsilon(\mathbf{x})} \exp\left(-\frac{\|\mathbf{x}-\mathbf{q}\|^2}{\epsilon^2}\right) e^{agg}}{\sum_{\mathbf{q} \in N_\epsilon(\mathbf{x})} \exp\left(-\frac{\|\mathbf{x}-\mathbf{q}\|^2}{\epsilon^2}\right)} \quad (3.3)$$

with

$$\epsilon = DT(\mathbf{x}) + \alpha. \quad (3.4)$$

Here, $N_\epsilon(\mathbf{x})$ contains all projections in D^m located within a radius ϵ from \mathbf{x} . We next draw $e^{agg}(\mathbf{x})$ as a RGBA texture, where the color components encode $e^{agg}(\mathbf{x})$ mapped via a suitable color map, and the transparency A is set to

$$A^{agg}(\mathbf{x}) = \begin{cases} 1 - \frac{DT(\mathbf{x})}{\alpha}, & \text{if } DT(\mathbf{x}) < \beta \\ 0, & \text{otherwise} \end{cases} \quad (3.5)$$

For $\alpha = 1, \beta = 1$, we obtain the classical colored scatterplot (Fig. 3.1 (a)). For $\alpha = 1, \beta > 1$, the space between projections is filled, up to a distance β , by the e^{agg} value of the closest data point. For $\alpha = 1, \beta = \infty$, we obtain a Voronoi diagram of the projections with cells colored by their e^{agg} values. This does not change the e^{agg} data values, but just displays them on larger spatial extents than individual pixels, making them easier to see, similarly to the use of Voronoi cells to show attributes in multidimensional projections in [24, 111]. This creates visualizations identical to those obtained by drawing scatterplots with point radii equal to β , without having the issues

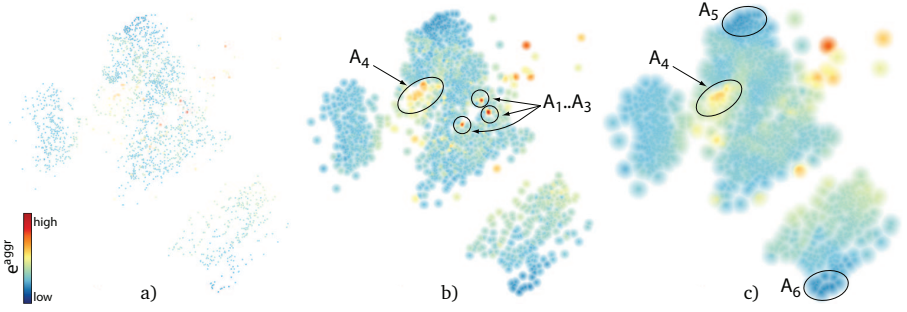


Figure 3.1: Aggregate error view, several levels of detail: (a) $\alpha = 1, \beta = 1$. (b) $\alpha = 5, \beta = 5$. (c) $\alpha = 20, \beta = 20$ pixels (see Sec. 3.2.2).

created by overlapping points. For $\alpha > 1, \beta > 1$, the result is similar to a smooth Shepard interpolation where the kernel size ϵ is given by the *local* point density. The parameter $\alpha \geq 0$ controls the *global* level-of-detail at which we visualize e^{agg} : Small values show more detail in dense point zones, but also emphasize small-scale signal variations that are less interesting. Larger α values create a smoother signal where coarse-scale error patterns are more easily visible.

Figs. 3.1 (b,c) show the aggregate error for the Segmentation dataset for various values of the parameters α and β . Here, $e_{ij} \in [-0.67, 0.35]$. The error range already tells that we have poorly projected points, but does not tell where these are. In Fig. 3.1 (b), with low values for both α and β , we see that e^{agg} is relatively smoothly distributed over the entire projection. However, we see three small red spots $A_1..A_3$. These are high-error outlier areas, which indicate points that are badly placed with respect to most other points. We also see a relatively high error area A_4 of larger spatial extent. Increasing both α and β produces a simplified visualization (Fig. 3.1 (c)). Larger β values fill in the gaps between points. Larger α values eliminate outlier regions whose spatial extent is smaller than α , such as the three small outlier areas $A_1..A_3$, but A_4 remains visible, since it is larger than α . We now also notice, better than in Fig. 3.1 (b), that the bottom and top areas (A_5, A_6) in the projection have dark blue values, with a significantly lower error than the rest of the projection.

Our image-based results are slightly reminiscent of the dense projection-precision-score (*pps*) maps of Schreck *et al.* [159]. Differences exist, however. First, our e_i^{agg} is a *global* metric, that tells how point i is placed with respect to all other points, whereas the *pps* metric characterizes *local* neighborhoods. Interpolation-wise, our technique (used with $\alpha = 1, \beta = \infty$) delivers the same Voronoi diagram as Schreck *et al.*, which is also identical to the space partitioning of the point-based Voronoi diagrams in [8, 111]. The data being mapped is, however, different: Our e^{agg} shows the sum of distance compression and stretching, whereas the techniques in [8, 111] treat these two quantities separately. In the next sections, we show how we split our aggregated insight into separate insights. Further on, both Schreck *et al.* and our method use smoothing to remove small-scale noise from such maps. However, whereas Schreck *et al.* uses a constant-radius smoothing kernel, which blurs the image equally strongly everywhere,

we use, as explained, a variable-radius kernel controlled by local density, which preserves better detail in non-uniform point clouds (scatterplots).

3.2.3 The False Neighbors view

While it is useful to assess the error distribution and find badly vs well-projected point groups, the aggregate error view does not tell us if the error is due to false neighbors, missing neighbors, or both. Let us first consider the false neighbors (case **A**, Sec. 3.1). To visualize these, we create a Delaunay triangulation of the projected point cloud that gives us the closest neighbors of each projected point in all directions, *i.e.*, the most important false-neighbor candidates for that point. To each edge E_k , $1 \leq k \leq 3$ of each triangle T of this triangulation, with vertices being the points \mathbf{q}_i and \mathbf{q}_j of D^m , we assign a weight $e_k^{false} = |\min(e_{ij}, 0)|$, *i.e.*, consider only errors created by false neighbors. Next, we interpolate e^{false} over all pixels \mathbf{x} of T by using

$$e^{false}(\mathbf{x}) = \frac{\sum_{1 \leq k \leq 3} \frac{1}{d(\mathbf{x}, E_k) \|E_k\|} e_k^{false}}{\sum_{1 \leq k \leq 3} \frac{1}{d(\mathbf{x}, E_k) \|E_k\|}} \quad (3.6)$$

where $d(\mathbf{x}, E)$ is the distance from \mathbf{x} to the edge E and $\|E\|$ is the length of the edge. Similarly to the aggregated error, we construct and render an image-based view for e^{false} as a RGBA texture. In contrast to the aggregated error, we use here a heated body colormap [74], with light hues showing low e^{false} values and dark hues showing high e^{false} values. This attracts the attention to the latter values, while pushing the former ones into the background. The transparency A is given by

$$A^{false}(\mathbf{x}) = A^{agg}(\mathbf{x}) \left(1 - \frac{1}{2} \left(\min \left(\frac{DT_T(\mathbf{x})}{DT_C(\mathbf{x})}, 1 \right) + \max \left(1 - \frac{DT_C(\mathbf{x})}{DT_T(\mathbf{x})}, 0 \right) \right) \right) \quad (3.7)$$

where $DT_T(\mathbf{x}) = \min(d(\mathbf{x}, E_1), d(\mathbf{x}, E_2), d(\mathbf{x}, E_3))$ is the distance transform of T at \mathbf{x} , $DT_C(\mathbf{x})$ is the distance from \mathbf{x} to the barycenter of T , and A^{agg} is given by Eqn. 3.5. The same technique is used in a different context to smoothly interpolate between two 2D nested shapes [153], to which we refer for further (simple) implementation details. The combined effect of Eqns. 3.6 and 3.7 is to slightly thicken, or smooth out, the rendering of the Delaunay triangulation. Note that this interpolation does *not* change the actual values e_k^{false} rendered on the triangulation edges. The distance-dependent transparency ensures that data is shown only close to the projection points.

Fig. 3.2 shows the false neighbors for the Segmentation dataset. Several aspects are apparent here. First, the rendering is similar to a blurred rendering of the Delaunay triangulation of the 2D projections colored by e^{false} , showing how each point relates to its immediate neighbors. Light-colored edges show true neighbors, while dark edges show false neighbors. Since edges are individually visible, due to the transparency modulation (Eqn. 3.7), we can see both the true and false neighbors of a point separately. The smooth transition between opaque points (on the Delaunay edges) and fully transparent points (at the triangles' barycenters) ensures that the resulting image is continuous and

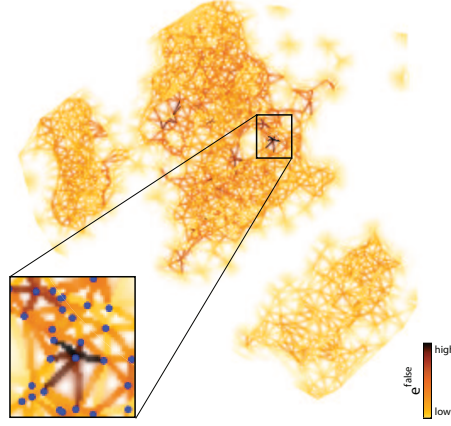


Figure 3.2: False neighbors view (see Sec. 3.2.3).

easier to follow at various screen resolutions than a Delaunay triangulation rendered with pixel-thin edges, as our edges appear slightly thicker.

In Fig. 3.2, two error-related aspects are visible. First, we see an overall trend from light to dark colors as we go further from the projection’s border towards the projection center. This confirms the known observation on DR methods that points projected near the border tend to be more accurate, since there is more freedom (and space) to place them. In contrast, projections falling deep inside the resulting point cloud tend to have more false neighbors, because the DR algorithm has there less space to shift points around to accommodate all existing distance constraints. Intuitively, we can think of this phenomenon as a ‘pressure’ which builds up within the projected point set from its border inwards. We shall see more examples of this phenomenon in Sec. 3.3. Secondly, we see a few small-scale dark outliers. Zooming in Fig. 3.2, we see that these are points connected by dark edges to most of their closest neighbors in a star-like pattern. Clearly, false neighbors exist here. These can be either the star ‘center’ or the tips of its branches. However, we also see that these tips have only one dark edge. Hence, they are too closely positioned to the star center only, and not to their other neighbors. Since the tip points are all positioned well with respect to their neighbors (except the star center), and the center point is positioned too closely with respect to all its direct neighbors, we can conclude that too little space was offered in the projection to the center point, or in other words that the center point is a false neighbor of its surrounding points.

The false neighbors view is related to Aupetit’s segment compression view, where the shortening of inter-point distances due to projection is visualized [8]. The underlying metrics, *i.e.* our e_{ij} (Eqn. 3.1) and m_{ij}^{distor} ([8], Sec. 3.2) are similar, up to different normalizations. However, the proposed visualizations are quite different. Aupetit uses so-called ‘segment Voronoi cells’ (SVCs). SVCs essentially achieve piecewise-constant (C^{-1}) interpolation of the values e_k^{false} , defined on the edges E_k of each Delaunay triangle T , over T ’s area, by splitting T in three sub-triangles using its barycenter. In contrast, our interpolation (Eqn. 3.6) is C^∞ over T . Also, our triangles are increasingly transparent

far away from their edges (Eqn. 3.7). Comparing our results (e.g. Figs. 3.2, 3.9 (a,d,g)) with SVCs (e.g. Figs. 7 (d), 12 (c) in [8]), we observe that SVCs exhibit several spurious elongated Voronoi cells that do not convey any information. Such cells do not exist in our visualization due to the transparency blending. Also, we argue that the artificial SVC edges linking projected points with Delaunay triangulation barycenters do not convey any information, but only make the visualization more complex. Such edges do not exist in our visualization due to our continuous interpolation.

3.2.4 The Missing Neighbors view

Besides false neighbors, projection errors (and subsequent misinterpretations) can also be caused by missing neighbors (case **B**, Sec. 3.1). Visualizing this by a space-filling method like for the aggregate error or false neighbors is, however, less easy. Given a projected point \mathbf{q} , its missing neighbors can be anywhere in the projection, and are actually by definition far away from \mathbf{q} . To locate such neighbors, we would need to visualize a many-to-many relation between far-away projected points.

We first address this goal by restraining the question's scope: Given a *single* point \mathbf{q}_i , show which of the other points $D^m \setminus \mathbf{q}_i$ are missing neighbors for \mathbf{q}_i . For this, we first let the user select \mathbf{q}_i by means of direct brushing in the visualization. Next, we compute the error $e_i^{\text{missing}} = \max_{j \neq i} (e_{ij}, 0)$, i.e., the degree to which \mathbf{q}_j is a missing neighbor for \mathbf{q}_i , and visualize e^{missing} by the same technique as for the aggregated error (Sec. 3.2.2).

Fig. 3.3 shows this for the Segmentation dataset, using the same heat colormap as in Fig. 3.2. In Figs 3.3 (a,b), we selected two points deep inside the central, respectively the lower-right point groups in the image. Since Figs. 3.3 (a,b) are nearly entirely light-colored, it means that these points have few missing neighbors. Hence, the 2D neighbors of the selected points are truly *all* the neighbors that these points have in nD . In Figs. 3.3 (c,d), we next select two points located close to the upper border of the large central group and the left border of the left group respectively. In contrast to Figs. 3.3 (a,b), we see now an increasingly darker color gradient as we go further from the selected points. This shows that points far away from these selections are actually projected *too* far, as they are actually more similar than the projection suggests. This is a known (but never visualized as such) issue of many DR methods, which have trouble in embedding high-dimensional manifolds in 2D: points close to the embedding's *border* are too far away from other points in the projection. Another interesting finding is that the color-coded Figs. 3.3 (c,d) do not show a *smooth* color gradient: We see, especially in Fig. 3.3 (c) that the colors appear grouped in several 'bands', separated by discontinuities. In other words, the projection method suddenly increases the error as we get over a certain maximal 2D distance.

The missing neighbors view is related to the proximity view of Aupetit [8]. In both views, a point i is selected and a scalar value, related to this selection, is plotted at all other points $j \neq i$. For Aupetit, this is the distance $m_j^{\text{prox}} = d^n(\mathbf{p}_i - \mathbf{p}_j)$ (normalized by its maximum). For us, it is the error e_j^{missing} . Both the distance and e^{missing} have, in general, the tendency to be small at points j close in 2D to the selected point i , and increase farther off from point i . However, the two quantities are different and serve different purposes. Visualizing m^{prox} is useful in finding points located within some

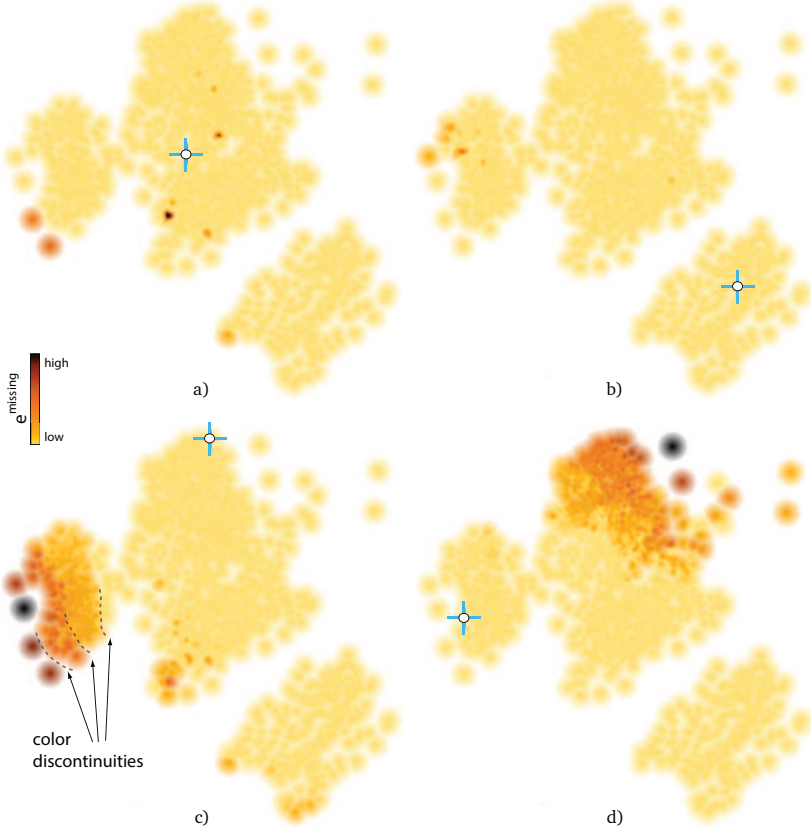


Figure 3.3: Missing neighbors view for different selected points. Selections are indicated by markers (see Sec. 3.2.4).

distance to the selection i . Finding projection errors is only *implicitly* supported, as these appear as non-monotonic variations in the m^{prox} signal. In contrast, $e^{missing}$ specifically emphasizes points projected too far, rather than conveying the absolute distance. Thus, our visualization helps locating projection errors rather than assessing proximity.

3.2.5 The Missing Neighbors Finder

Although providing details for single points, the views in Sec. 3.2.4 cannot show missing neighbors for an entire dataset. We address this goal by a different method, as follows. Consider all positive values of e_{ij} . By definition, these give all point-pairs which are projected too far away. We sort these values decreasingly, and select the largest ϕ percent of them, where ϕ is a user-provided value. The selected values give the point pairs which are worst placed in terms of overestimating their true similarity. We next construct a graph $G = (V, E)$ whose nodes V are the projected points \mathbf{q}_i present in such point pairs, and edges E indicate the pairs, with e_{ij} added as edge weights. Next, we draw

G using the KDEEB edge bundling technique [85], which provides robust, easy to use, and real-time bundling of graphs with tens of thousands of edges on a modern GPU. We color the bundled edges based on their weight using a grayscale colormap (with white mapping low and black mapping high weights), and draw them sorted back-to-front on weight and with an opacity proportional to the same weight. The most important edges thus appear atop and opaque, and the least important ones are at the bottom and transparent.

Fig. 3.4 shows this visualization, which we call the *missing neighbors finder*, with bundles that connect a single selected point with its most important missing neighbors (bundles connecting multiple points are discussed later on). The background images show $e^{missing}$ (Sec. 3.2.4). Dark bundle edges attract attention to the most important missing neighbors. For the selected points in images (a) and (b), we see that there are only very few and unimportant missing neighbors (few half-transparent edges). For the selected points in images (c) and (d), the situation is different, as the bundles are thicker and darker. Bundle fanning shows the *spread* of missing neighbors for the selected points: In image (c), these are found mainly in the left point group, with a few also present in the lower part of the central group. In contrast, all missing neighbors of the point selected in image (d) are at the top of the central group.

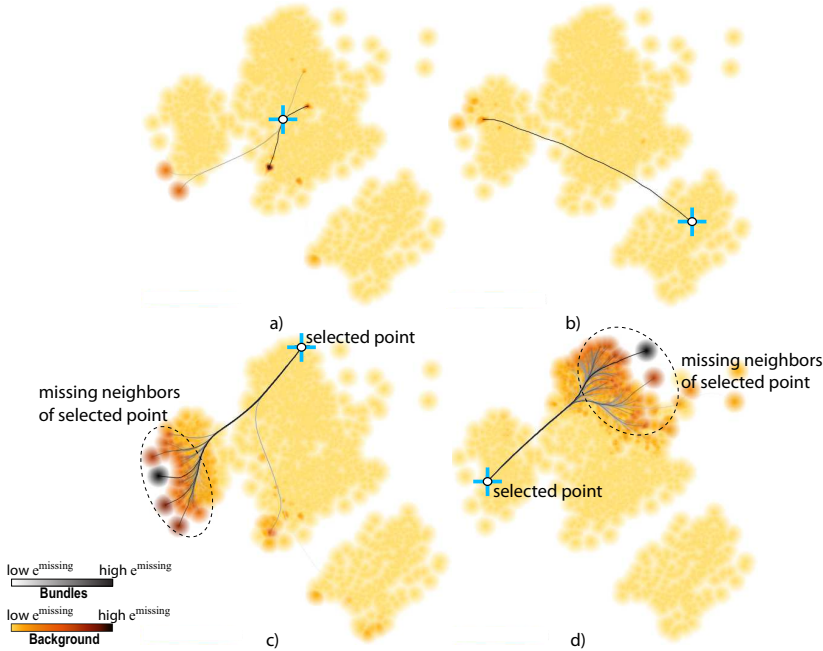


Figure 3.4: Missing neighbors finder view for four selected points. Selections are indicated by markers (see Sec. 3.2.5).

The main added value of the missing neighbors finder appears when we visualize the many-to-many relations given by all projected points. Fig. 3.5 shows this result for three

values of ϕ for the Segmentation dataset. The background shows now the aggregated error (e^{agg} , Sec. 3.2.2). We color bundles from black for largest error e_{ij} to white for largest error above the user-provided parameter ϕ . Image (a) shows the $\phi = 1\%$ worst missing-neighbor point-pairs. These link the top-right area of the central group with the left frontier of the left group. Adding more missing neighbor pairs to the view (image (b), $\phi = 3\%$) strengthens this impression. Adding even more missing neighbor pairs (image (c), $\phi = 20\%$) reveals additional missing-neighbor pairs between the two areas indicated above (light gray parts of thick top bundle), and also brings in a few missing neighbors between these areas and the lower-right point group (light gray thin bundle going to this group). Nearly all bundles appear to connect point pairs located on the *borders* of the projection. This strengthens our hypothesis that such point pairs are challenging for the LAMP projection, which we noticed using the interactive missing neighbors view (Sec. 3.2.4). However, as compared to that view, the bundled view shows all such point pairs in a single go, without requiring user interaction.

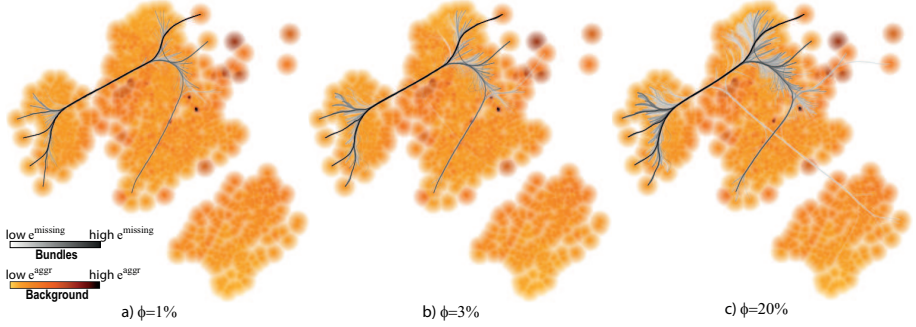


Figure 3.5: Missing neighbors finder view, all point pairs, for different ϕ values (see Sec. 3.2.5).

3.2.6 The Group Analysis views

As outlined in Sec. 3.1, the false and missing neighbors issues for individual points become, at group level, the problems of false and missing group members respectively. We next propose two visualizations that assist in finding such issues.

First, let us refine the notion of a group. Given the tasks in Sec. 3.1 (C), a group $\Gamma \subset D^m$ is a set of projected points which form a *visually* well-separated entity. When users see points in a group, they understand that these share some commonality, but are different from points in other groups. In the LAMP projection of our Segmentation dataset, we see three such groups (Figs. 3.1-3.5). Group perception is, obviously, subject to many factors such as user preferences and level-of-detail at which one focuses. However, once a user has established which are the groups (s)he sees in a visualization, the false and missing membership issues become relevant.

We allow users to select groups in a given projection by several mechanisms: direct interactive selection, mean-shift clustering [36], and upper thresholding of the point density [52]. Other user-controlled methods can be used if desired, e.g., K -means

or hierarchical agglomerative clustering, *e.g.* [92, 91]. The actual group selection mechanism is further of no importance to our visualization method. We next render each obtained group $\Gamma = \{\mathbf{q}_i\}$ by the shaded cushion technique in [182] as follows. First, we compute a density map $\rho(\mathbf{x}) = \sum_{\mathbf{q} \in \Gamma} K(\mathbf{x} - \mathbf{y})$, where K is an Epanechnikov kernel of width equal to the average inter-point distance δ in Γ , following [36]. Next, we compute a threshold-set Γ_δ of ρ at level δ , and its distance transform DT_{Γ_δ} . Finally, we render a RGBA texture over Γ_δ , where we set the color a fixed hue (light blue in our case) and the transparency A to $\sqrt{DT_{\Gamma_\delta}}$ (following the approach in [182] which shows that such a transparency profile creates naturally-looking cushions).

Having now groups both as a data structure and also shown in the visualization, we adapt the missing neighbors and finder techniques (Secs. 3.2.4, 3.2.5) to show missing group members. For this, we compute a value

$$e_\Gamma^{\text{missing}}(\mathbf{q}_i) = \begin{cases} \min_{\mathbf{q}_j \in \Gamma} (e_{ij}) & \text{if } \mathbf{q}_i \notin \Gamma \\ 0 & \text{otherwise} \end{cases} \quad (3.8)$$

at each projected point \mathbf{q}_i , and visualize $e_\Gamma^{\text{missing}}$ using the same technique as for missing neighbors.

Fig. 3.6 (a,b) show two missing group members views. The shaded cushions show the three groups identified in our Segmentation dataset. Several points fall outside of all groups. This is normal, in general, *e.g.* when the user cannot decide to which group to associate a point. In image (a), we select the bottom group Γ_{bottom} . The underlying color map shows now $e_{\Gamma_{\text{bottom}}}^{\text{missing}}$ (Eqn. 3.8). All points appear light yellow. This means that, with respect to Γ_{bottom} *seen as a whole*, no points are projected too far, so Γ_{bottom} has no missing members. In image (b), we do the same for the left group Γ_{left} . The image now appears overall light yellow, except for a small dark-red spot in the upper-right corner of the central group Γ_{center} . Here are a few points which are placed too far from any point in Γ_{left} . These are highly likely to be missing members of Γ_{left} . To obtain more insight, we now use the bundle view in Sec. 3.2.5, with two changes. First, we build only bundles that have an endpoint in the selected group. Secondly, we consider all edges rather than showing only the most important ones. Image (c) shows the bundle view for Γ_{bottom} . We see only a few bundled edges, ending at a small subset of the points in Γ_{bottom} . This strengthens our hypothesis that there are no points outside Γ_{bottom} which should be placed closer to *all* points in Γ_{bottom} – or, in other words, that Γ_{bottom} has no missing members. Image (d) shows the bundled view for Γ_{left} . The bundle structure tells us that the top-right part of Γ_{center} contains many missing neighbors of Γ_{left} . In particular, we see dark bundle edges that connect to dark-red points. This is a strong indication that these points can indeed be missing members of Γ_{left} . For a final assessment, the user can interactively query the discovered points' details (attribute values) and, depending on these, finally decide if these points are missing group members or not.

3.2.7 The Projection Comparison view

Consider running the same DR algorithm with two different parameter sets, or projecting a dataset by two different DR algorithms. How to compare the results from the viewpoint

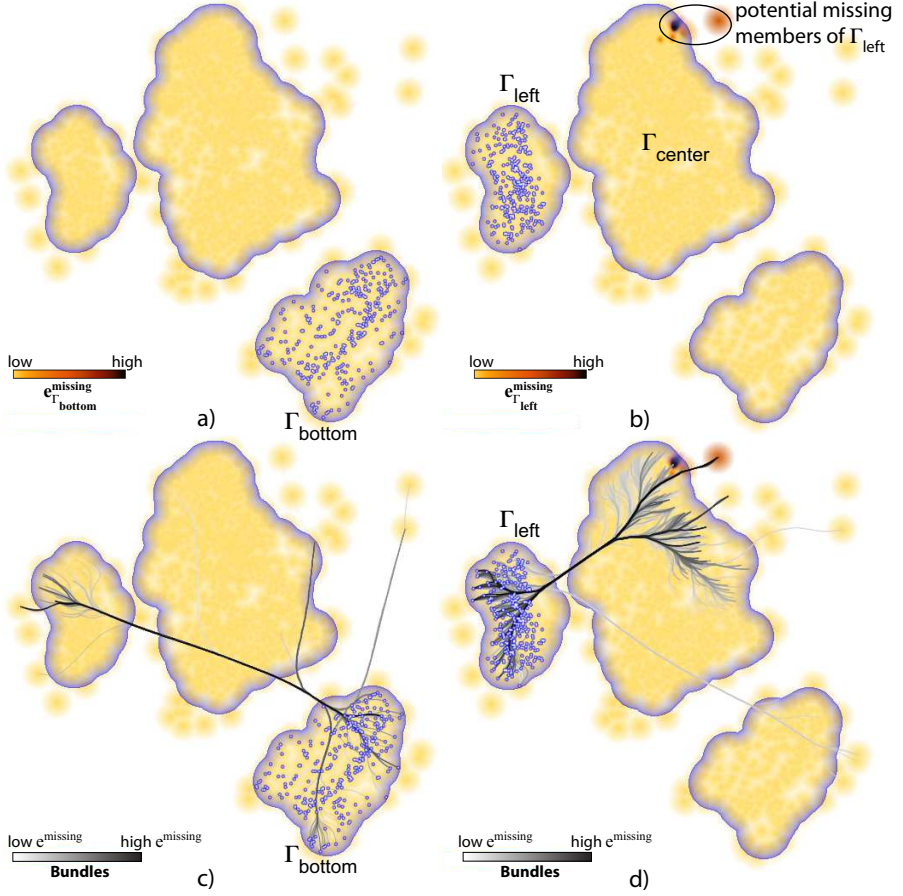


Figure 3.6: Missing members for two point groups. Points in the selected groups are drawn as marked (see Sec. 3.2.6).

of distance preservation? Subsequent questions are: Which points that were (correctly) placed close to each other in one projection are now ‘pulled apart’ in the other projection? Do the two projections deliver the same groups of points? Understanding all these aspects is crucial to further using the respective projections for interpreting the original high-dimensional data. Indeed, if we were able to see that specific point-groups are placed at significantly different locations in the 2D space by different projection techniques and/or parameter settings, this would mean that great care should be invested in assessing the respective tools (projection techniques and/or parameter settings) *prior to* interpreting the projection results. If, in contrast, the main point-groups in the projection would be the same with respect to technique and/or parameter settings, this would mean that one could use the respective projection tools very much like agnostic ‘black boxes’ further on.

To answer such questions, we propose the *projection comparison* view. The view reads two projections D_1^m and D_2^m of the same input dataset D^n . For each point-pair

$(\mathbf{q}_i^1 \in D_1^m, \mathbf{q}_i^2 \in D_2^m)$, we compute a displacement

$$e_i^{disp} = \frac{\|\mathbf{q}_i^1 - \mathbf{q}_i^2\|}{\max_i \|\mathbf{q}_i^1 - \mathbf{q}_i^2\|}. \quad (3.9)$$

We next build a graph whose nodes are points in $D_1^m \cup D_2^m$. Edges relate point pairs $(\mathbf{q}_i^1 \in D_1^m, \mathbf{q}_i^2 \in D_2^m)$, and have the values e^{disp} as weights. We visualize this graph via edge bundling, as for the missing neighbors finder (Sec. 3.2.5).

Fig. 3.7 (a) shows a view where we compare the Segmentation dataset projected via LAMP (red points, D_1^m) and LSP (green points, D_2^m). The two projections are quite similar, since red and green points occur together in most cases. However, this image does not tell if the two projections create the same *groups* of points, since we do not know how red points match the green ones. Fig. 3.7 (b) shows the projection comparison view for this case. We immediately see a thin dark bundle in the center: This links corresponding points which differ the most in the two projections. Correlating this with image (a), we see that LSP decided to place the respective points at the bottom (A_{LSP}) of the central group, while LAMP moved *and* also spread out these points to the top (A_{LAMP}). However, points around the locations A_{LSP} and A_{LAMP} do not move much between the two projections, as we see only light-colored bundles around these locations, apart from the dark bundle already discussed. Hence, the motion of these points indicates a neighborhood problem in one or both of the projections. Indeed, if e.g. the points in A were correctly placed by LAMP (into A_{LAMP}), then the decision of LSP to move the point-group A all the way up in the visualization (to A_{LSP}) should also have moved the *neighbors* of A_{LAMP} . Since this does not happen, A_{LSP} cannot be close to the same points that A_{LAMP} was. A similar reasoning applies if we consider that A_{LSP} is correct – it then follows that A_{LAMP} cannot be correctly placed with respect to its neighbors.

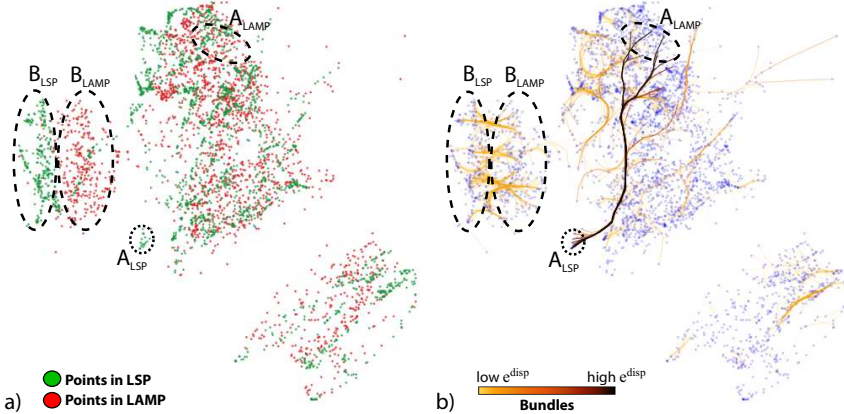


Figure 3.7: Comparison of two projections. (a) LAMP (blue) and LSP (red) points. (b) Bundles show corresponding point groups in the two projections (see Sec. 3.2.7).

Apart from this salient dark-colored bundle, we see many shorter and light-colored

bundles. These show smaller-scale displacements between the two projections. For instance, we see how the red points at the right of the left group (B_{LAMP}) are moved to the left (B_{LSP}) of the same group. As these bundles fan out relatively little, do not have many crossings, and they are short, it means that B_{LSP} is almost a *translation* to the left of B_{LAMP} , so the two projections depict the same structure of the left group. Also, we do not see any bundle exiting this left group. This means that both LAMP and LSP keep all points in this group together. Finally, in the bottom-right group we see just a very few short light-colored bundles. Most points in this group do not have any bundles connected to them. This means that e_i^{disp} for these points is very small (yielding thus very short, nearly transparent, bundles). From this, we infer that LAMP and LSP produce very similar layouts for this group. If users are interested only to spot the most salient differences between two projections, and want to ignore such small-scale changes, this can be easily obtained by mapping e_i^{disp} to bundle-edge transparency.

3.2.8 Usage scenario

Considering that the user is offered quite a few different views to analyse projection errors, each with specific features and goals, the next question arises: How to put all these views together to form a coherent *usage scenario* for a common analysis task? Below we propose such a usage scenario. The view names herein refer to the respective techniques presented earlier in this section.

Step 1: Start with the *Aggregated Error* view. This shows an overview of the error at all points, without a distinction between false or missing neighbors. Next, check if (a) there are regions or groups with substantial errors or (b) the overall error is low. Case (b) indicates that the projection is quite good and that nothing else needs to be improved. In case (a), continue with steps 2, 3, and 4.

Step 2: The *Missing Neighbors Finder* view can be enabled and disabled freely over the *Aggregated Error* view to show the most important missing neighbors between all points. The user should notice now whether this view shows bundles having high error values (*i.e.* dark-colored). If so, there are important missing neighbors between the groups connected by such bundles. These groups must be further analysed with the *Group Analysis Views*. If not, *i.e.* the bundles are colored (light) gray, this tells that the projection is good and, although there are missing neighbors, they are in a low error range and should not threaten the projection interpretation.

Step 3: Points, groups or regions found problematic in steps 1 and 2 are now analysed in more detail using the *False Neighbors* and *Missing Neighbors* views. For groups detected in step 1 the most important thing is to find out exactly what kind of error is present: Are they (a) wrongly placed with respect to each other and other close points (false neighbors) or (b) in relation to far away points that should be closer (missing neighbors)? For groups detected in step 2, the error is already identified from the beginning: They have a high rate of missing neighbors. In this case, the question to be answered is: Which

points are exactly the problematic ones inside the detected groups, or where exactly do the relations (bundle edges) with the highest errors start and end from? By using these two views, the user should be able to establish exactly which are the more problematic points (or groups), and what kind of error these have.

Step 4: Knowing now where exactly errors occur, we consider the next questions: (1) Are such errors really a problem? (2) Do they show unexpected results related to how the projection should work with the provided data? (3) Are the problematic points important for the analysis task at hand? If questions (1-3) all answer ‘no’, then we have a good projection for our data and analysis task, and our analysis stops. If any question (1-3) answers yes, then the user must improve the projection of problematic points, as follows. If the user is a projection designer testing the accuracy of a new method, (s)he should go back to the algorithm and use the new insight gotten from this analysis to improve that algorithm. If the user has no access to the projection implementation, the solution is to re-execute the analysis from step 1 with either (i) a new projection algorithm that might better fit the specific data and task; or (ii) a new set of parameters for the same algorithm. The new results can be compared with the old ones to determine if the errors have decreased or if the errors moved into a new region where they are not as important for the task at hand. For the second task, the *Projection Comparison View* can be used.

3.3 Applications

We now use our views to study several projections for several parameter settings – thus, to explore the space P that controls the creation of a DR projection. First, we present the datasets used (Sec. 3.3.1), the studied projection algorithms (Sec. 3.3.2), and their parameters (Sec. 3.3.3). Next, we use our views to explore the considered parameter settings (Secs. 3.3.4, 3.3.5).

3.3.1 Description of datasets

Apart from the Segmentation dataset used so far, we consider the following datasets:

Freephoto: contains 3462 images grouped into 9 unbalanced classes [58]. For each image, we extract 130 BIC (border-interior pixel classification) features. Such features are widely used in image classification tasks [173].

Corel: composed of 1000 photographs that cover 10 specific subjects. Similarly to the Freephoto dataset, we extract for each image a vector of 150 SIFT descriptors [115].

News: contains 1771 RSS news feeds from BBC, CNN, Reuters and Associated Press, collected between June and July 2011. The 3731 dimensions were created by removing stopwords, employing stemming and using term-frequency-inverse-document-frequency counts. We manually classified the data points based on the perceived main topic of the news feed resulting in 23 labels. Given the imprecision of the manual classification and

the restriction to have one topic per point, the labels are unbalanced for a number of points. Also, for other points (with different labels), we can still have a high similarity of content.

Sourceforge: This publicly available dataset contains 24 software metrics computed on 6773 open-source C++ software projects from the sourceforge.net website [121]. Metrics include classical objet-oriented quality indicators such as coupling, cohesion, inheritance depth, size, complexity, and comment density [107], averaged for all source code files within a project.

3.3.2 Description of projections

We detail next the projection algorithms whose parameter spaces we will next study. We chose these particular algorithms based on their availability of documented parameters, scalability, genericity, presence in the literature, and last but not least availability of a good implementation.

LSP: The Least Squares Projection [139] uses a force-based scheme to first position a subset of the input points, called control points. The remaining points in the neighborhood of the control points are positioned using a local Laplace-like operator. Overall, LSP creates a large linear system that is strong in local feature definition. LSP is very precise in preserving neighborhoods from the nD space to the 2D space.

PLMP: The Part-Linear Multidimensional Projection (PLMP) [141] addresses computational scalability for large datasets by first constructing a linear mapping of the control points using the initially force-placed control points. Next, this linear mapping is used to place the remaining points, by a simple and fast matrix multiplication of the feature matrix with the linear mapping matrix.

LAMP: Aiming to allow more user control over the final layout, the Local Affine Multidimensional Projection (LAMP) [95] provides a user-controlled redefinition of the mapping matrix over a first mapping of control points. LAMP also works by defining control points, which are used to build a family of orthogonal affine mappings, one for each point to project. LAMP has restrictions regarding the number of dimensions against the number of points. Also, LAMP cannot directly work with distance relations, *i.e.*, it needs to access the nD point coordinates. However, LAMP is very fast, without compromising the precision reached, for instance, by LSP. Both LSP and LAMP can be controlled by a number of parameters, such as the control point set.

Pekalska: Another class of projection techniques works with optimization strategies. These are, in general, quite expensive computationally. To improve speed, Pekalska *et al.* [144] first embeds a subset of points in 2D by optimizing a stress function. Remaining points are placed using a global linear mapping, much like LAMP and LSP.

ISOMAP: The ISOMAP technique [184] is an extension of classical Multidimensional Scaling (MDS) that aims to capture nonlinear relationships in the dataset. ISOMAP replaces the input distance between point pairs by an approximation of the geodesic distance given by the shortest path on a graph created connecting neighbor points in the original space with the original distance as weight. The final 2D coordinates are computed via a conventional MDS embedding with calculations of eigenvalues over the distance relations of the previous step.

3.3.3 Description of parameters to analyse

Most techniques that initially project control points use a simplified iterative force-based algorithm, such as the one of Tejada *et al.* [179]. The number of *iterations* of force-based placement influences the control points' positions, and is, thus, a relevant parameter. LSP control points are typically the centroids of clusters obtained from a clustering of the input dataset. The *number of control points* is thus a second relevant parameter for LSP. To position points in the neighborhood of a given control point, LSP solves a linear system for that neighborhood. The neighborhood size (*number of neighbors*) is a third relevant parameter.

In LAMP, the affine mappings are built from a neighborhood of control points. The size of the control point set used to build the mapping, expressed as a *percentage* of the size of the control point set, is the main parameter here. The choice of control points and the choice of the initial projection of the control points are also parameterizable, just as for LSP, PLMP, and Pekalska. However, in LAMP, these parameters are mainly interactively controlled by the user, and thus of a lesser interest to our analysis.

ISOMAP, just as the previous methods, also requires the expression of neighborhoods. The main, and frequently only, exposed parameter of ISOMAP is the number of *nearest neighbors* that defines a neighborhood.

3.3.4 Overview comparison of algorithms

To form an impression about how the goals outlined in Sec. 3.1 are better, or less well, satisfied by LAMP, LSP, PLMP, and Pekalska, we start with an overview comparison.

Figure 3.8 shows the false neighbors, aggregated error, and most important $\phi = 5\%$ missing neighbors for the Segmentation dataset. To ease comparison, color mapping is normalized so that the same colors indicate the same absolute values in corresponding views. The aggregate error (top row) is quite similar in both absolute values and spread for all projections, *i.e.*, lower at the plot borders and higher inside, with a few dark (maximum) islands indicating the worse-placed points. Overall, thus, all studied projections are quite similar in terms of distance preservation quality. The false neighbors views (middle row) show a similar insight: Border points have few false neighbors (light colors), and the density of false neighbors increases gradually towards the projections' centers. Although local variations exist, these are quite small, meaning that all studied projections are equally good from the perspective of (not) creating false neighbors.

The missing neighbors view (bottom row) is however quite different: By looking at the size and color of the depicted bundles, we see that LSP and Pekalska have many more

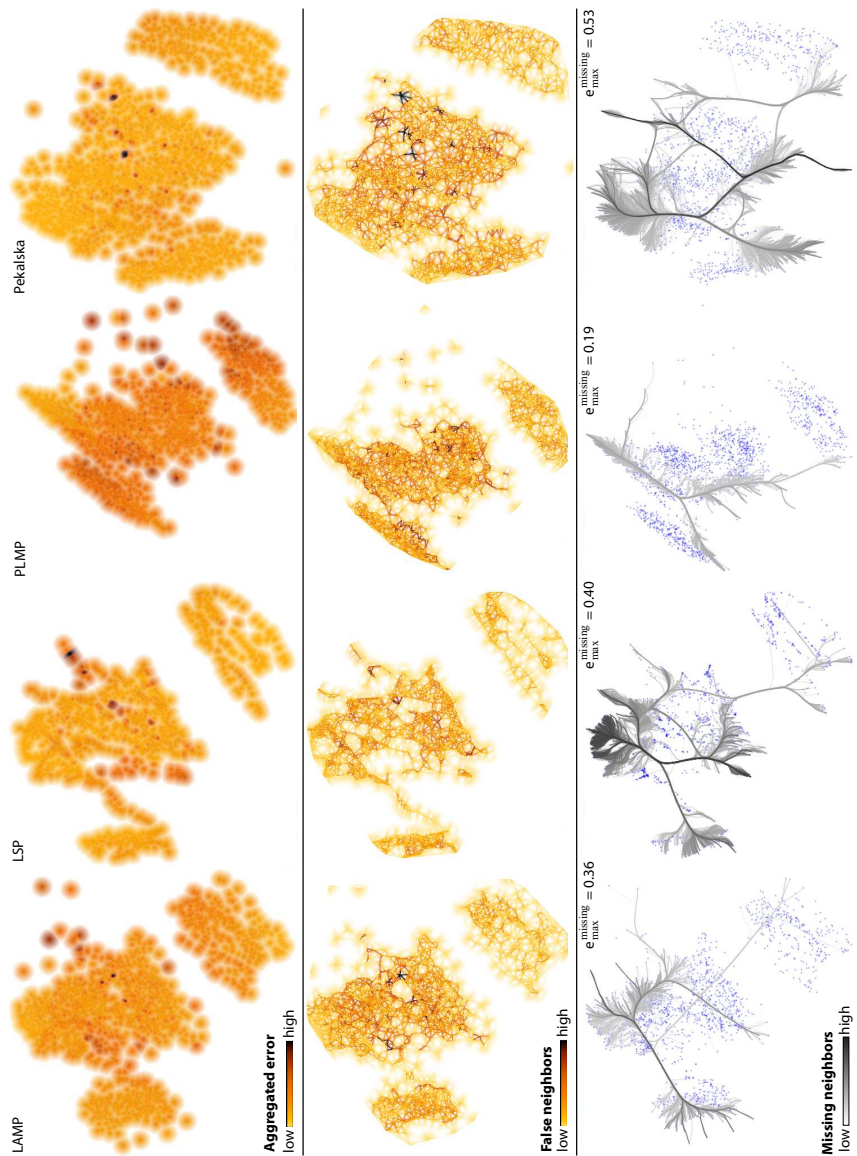


Figure 3.8: Comparison of LAMP, LSP, PLMP, and Pekalska projections for the Segmentation dataset (see Sec. 3.3.4)

important missing neighbors than PLMP, while LAMP has the fewest missing neighbors. In all cases, we see bundles that connect borders of the projected point-set. This confirms that all studied projections optimize placement of close points than far-away points. We also see that the missing neighbors are spread differently over the data: For LAMP, there are no bundles going to the bottom-right point cluster, showing that this cluster is indeed well separated in the projection, as it should be in relation to the n D data. In contrast, LSP, PLMP, and Pekalska all have bundles going to this cluster, indicating that they place these points too close to the remaining projected points.

3.3.5 Parameter analysis

We next refine our overview analysis by selecting two of the studied algorithms: LAMP and LSP. We next vary several of their parameters, and evaluate the resulting projections' quality with respect to this variation.

LAMP – Different control point percentages: Fig. 3.9 shows the results of LAMP for the Freephoto dataset with three different values for the *percentage* parameter: 10%, 30% and 50%. The error has been normalized on each view type (column in the figure).

First, we see that the final layout of the point cloud does not change drastically while varying the *percentage* parameter, only showing a 90 degree clockwise rotation for the value of 30%. While analysing the false neighbors view, we also see that, while the light brown areas are large – meaning that a moderate amount of error can be expected on the whole layout – the dark-colored spots are found nearer to the center. This suggests that LAMP positions the most problematic points in the center, surrounded by the rest of the points. By focusing on the dark spots (points with the largest false neighbor errors) throughout the parameter variation we can see that the value of the largest errors on each result remain similar – no view has many more, or much darker-colored, areas.

For the missing neighbors view, we selected a point near the upper border of the layout, marked by a cross in Figs. 3.9 (b), (e) and (h)), since missing neighbors occur mainly on the borders of the projection, as we have already observed in Section. 3.2.4. The dark spot in Fig. 3.9 (h) is where the largest error occurs over these three views. While in Fig. 3.9 (b) there are a few orange spots showing moderate error, in Fig. 3.9 (e) the error decreases considerably, and then increases again in Fig. 3.9 (h). This suggests that using about 30% of neighbors is a good value for avoiding large numbers of missing neighbors. We confirmed this hypothesis on several other datasets (not shown here for brevity). Finally, the aggregated error view shows results very similar to the false neighbors view: More problematic points (dark spots) are pushed to the center, and moderate error is found spread evenly over the entire layout. This shows that, for LAMP, most errors come from false neighbors rather than from missing neighbors.

LSP – Different numbers of control points: Figure 3.10 shows the same dataset (Freephoto) projected with LSP. The varying parameter is the *number of control points*. We use here the same views as in Fig. 3.9, and normalized the error in each column. By looking at the false neighbors views, we see a spatial interleaving of light-yellow and orange-brown colored areas in the projection. This contrasts with LAMP (Fig. 3.9) where

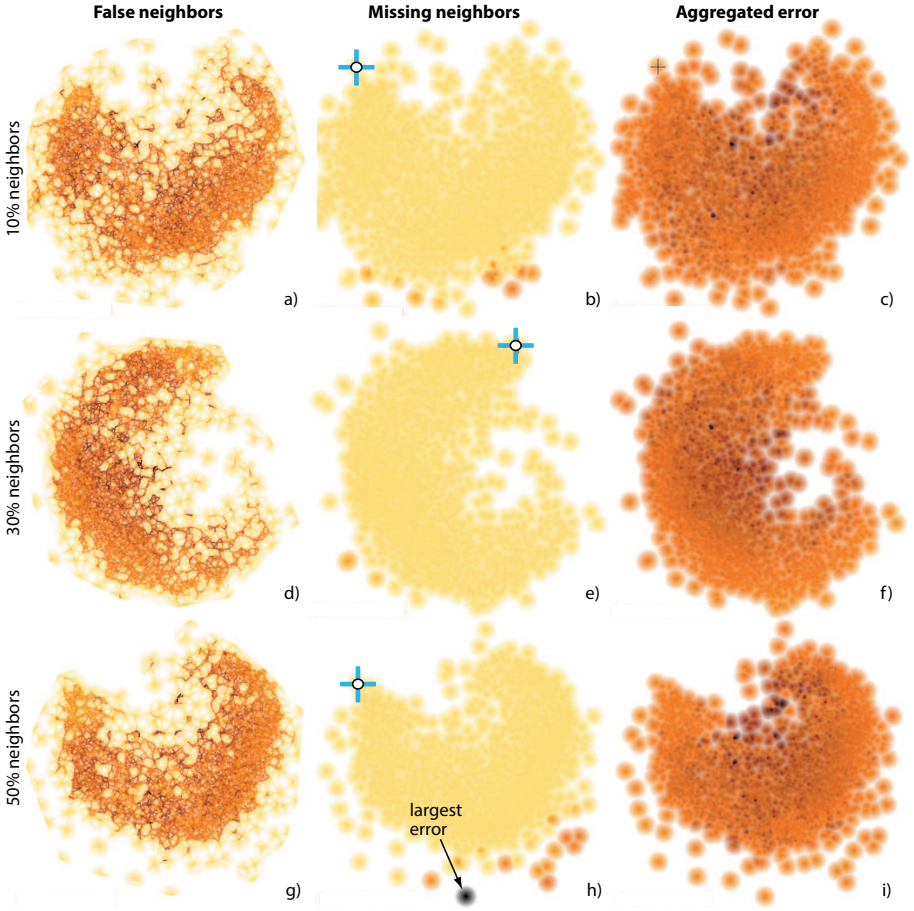


Figure 3.9: Applications – LAMP algorithm, Freephoto dataset, different neighbor *percentages* per row (see also Fig. 3.10).

the larger missing neighbor errors are consistently located away from the projection border. As the *number of control points* increases, the large error areas get more compact and closer to the projection center, but we see no increase in error severity (the amount of the orange and dark-red spots stays the same). In the missing neighbors views, the dark-colored areas in Fig. 3.10 (b) disappear largely in images (e) and (h), which means that the missing neighbors severity decreases when our control parameter increases. Comparing this with LAMP (Fig. 3.9 b,e,h), this shows that LAMP and LSP behave in opposite ways when dealing with missing neighbors. Finally, like for LAMP, the aggregate error views show the worst errors (dark spots) located in the center: The most problematic points are pushed inside by the other points which surround them, creating a mix of both false neighbors and missing neighbors. The severity of the errors, however, does not change visibly between the three parameter values.

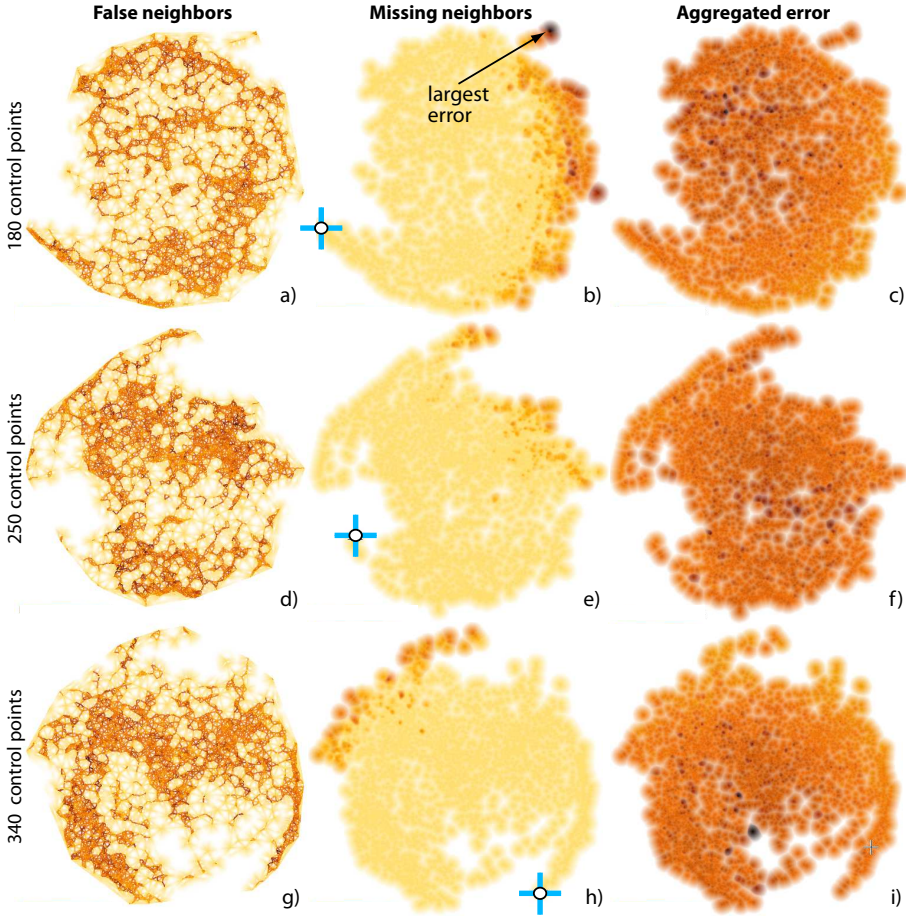


Figure 3.10: Applications – LSP technique, Freephoto dataset, different numbers of *control points* per row (compare with Fig. 3.9)

LSP – Different numbers of neighbors: We next examine the effect of a second parameter of LSP: *number of neighbors*. For the Freephoto dataset, we fix 250 control points and vary the number of neighbors to 10, 50 and 100. Fig. 3.11 shows the results with the missing neighbors finder view. We see that the most significant errors are initially concentrated between groups A, B and C, with C being essentially too far placed from both A and B. Increasing our parameter reduces has a positive impact on solving the missing neighbors problem between groups A and C, bringing them together into the group marked AC. The main missing neighbors are now concentrated in the relationship between groups AC and B. The ‘concentration’ of error given by the parameter increase is, upon further analysis, explainable by the working of LSP: Given a neighborhood N , LSP’s Laplace technique positions all points in N close to each other in the final layout.

However, the position of the neighborhoods N_i *themselves* is given only by the control points, which are determined by the initial force-based layout. If this layout suboptimally places two control points i and j too far away from each other, then *all* points within the neighborhoods N_i and N_j end up being too far away from each other. Hence, as the neighborhood size increases, the likelihood to see fewer thick high-error bundles increases. This insight we found is interesting since it was not reported in the LSP literature so far, and it can be explained (once we are aware of it) by the algorithmics of LSP.

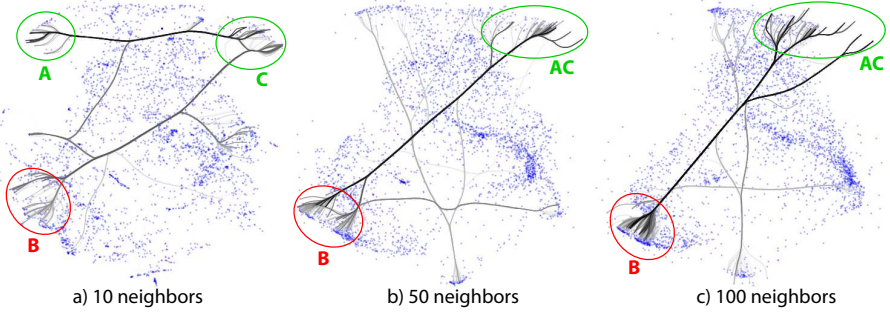


Figure 3.11: Applications – LSP technique, Freephoto dataset, different numbers of *neighbors*. Bundles show most important missing neighbors.

LAMP – Different datasets: We next analyse the LAMP technique applied to three different datasets: Corel (1000 elements), Freephoto (3462 elements), and Sourceforge (6773 elements). The varying parameter is now the input *dataset* itself. The aim is to see whether (and how) errors are affected by the nature of the input data, *e.g.* distribution of similarity, number of dimensions, and number of points. Figure 3.12 top row shows the false neighbors views. We see here that, while for the first two datasets the behavior of false neighbors is similar to earlier results, for the largest dataset (Sourceforge) there are much fewer false neighbors. These are located close to the intersection area of the two apparent groups in the image, and on the borders of these groups. This, and the low errors (light colors) inside the groups may indicate that both groups have a high degree of cohesion between their inner elements. The large errors on close to the intersection areas and borders can indicate elements that could be in either group, respectively very different from all other elements. Figure 3.12 (a) shows a similar pattern: Most false neighbors are located at the ‘star’ shape’s center, while the arms of the star contain elements that are more cohesive. This may indicate that the dataset contains a number of cohesive groups equal to the number of star arms, and elements in the center belong equally to all groups.

While analysing the missing neighbors for several points selected on the periphery of the projections, we see that the errors are smaller for Figs. 3.12 (d) and (e), and considerably larger for Fig. 3.12 (f). For the last image, we selected a point close to the intersection area of the perceived groups. Image (f) shows that this point is *equally* too far placed from most points in *both* perceived clusters. The size and speed of increase of

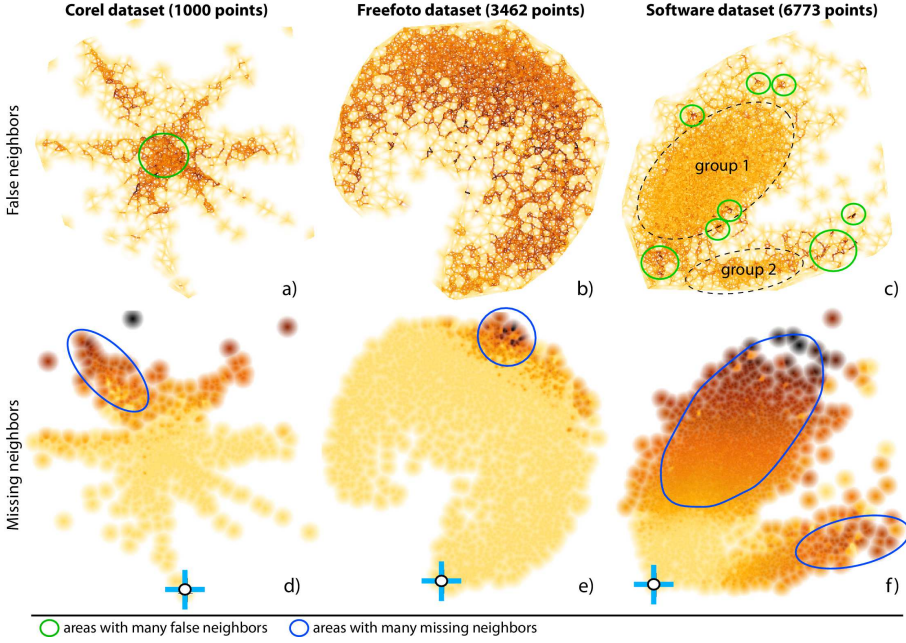


Figure 3.12: Applications – One algorithm (LAMP), different datasets. Top row: false neighbors. Bottom row: missing neighbors.

the error (as we get further from this point in the projection space) strongly suggests that the selected point belongs stronger to both perceived groups than the projection indicates. This strengthens our initial hypothesis that the area separating the two groups belongs equally to these groups.

ISOMAP – Different numbers of neighbors: To illustrate a different type of analysis made possible by our work, Fig. 3.13 shows the effect of changing the number of *neighbors* in ISOMAP on missing group members. Our group Γ of interest, shown first on Fig. 3.13 (a), is highlighted in images (b-d) by a shaded cushion. Besides the fact that Γ moves from the left of the projection to the right, images (b-d) show how its missing members behave as we change our parameter. At first, in Fig. 3.13 (b), we see that the most important missing neighbors are found in two other areas A_1 and A_2 on the far side of the layout. We also notice many black edges, which means that the points in A_1 and A_2 are indeed too far away from all points in the selected group. The relatively large fan-out of the bundles show that the group misses many members, and these are scattered widely over the projection. As the parameter increases, we see in image (c) that the missing members spread out even more, but the severity of the errors decreases (as shown by the lighter colors of $e_{\Gamma}^{missing}$ background). The inner fanning of the edges, inside Γ , is still large, which shows that many group members miss neighbors. Finally, in Fig. 3.13 (d), issues decrease significantly: We see thinner bundles, which imply less error; the bundle fanning inside Γ is relatively small, meaning that most of

Γ 's points do not miss neighbors; and the fan-out of the bundles is smaller, showing that the missing group members are now more concentrated than for the first two parameter values. This leads to the conclusion that, for the analysed group, the increase of the number of neighbors parameter has a positive impact on the final projection quality.

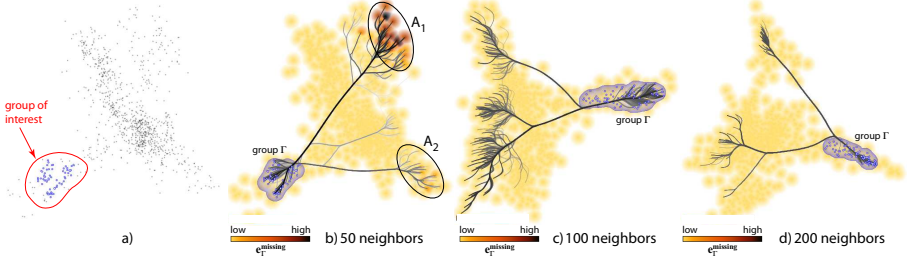


Figure 3.13: Applications – ISOMAP projection, finding missing group members for different numbers of *neighbors*.

LSP – Different numbers of iterations: The final analysis we present compares two different LSP projections of the same dataset (News), computed using values of 50, respectively 100 for the *number of iterations* parameter of the control-point force-directed placement.

Figures 3.14 (a) and (b) show the two LSP projections. In each of them, several high-density groups are visible. These are strongly related news feeds, *i.e.*, which likely share the same topic (see Sec. 3.3.1). However, without extra help, we cannot *relate* the two projections, *e.g.*, find out (a) if points significantly change places due to the parameter change; (b) which groups in one projection map to groups in the other projection; and (c) whether points in a group in one projection are also grouped in the second projection.

To answer question (a), we use the projection comparison view (Sec. 3.2.7). The result (Fig. 3.14 (c)) shows that there are many large point shifts; the bundle criss-crossing also shows that groups change places in the projection. This is a first indication that LSP is not visually stable with respect to its number of iterations parameter. Next, we manually select three of the most apparent point groups in one projection, shown in Fig. 3.14 (a) by the shaded cushions A, B, C . We examine these in turn. In Fig. 3.14 (d), we show how points in group A shifted, in the second projection, to a group A_1 . Virtually all bundled edges exiting A end in A_1 , so the parameter change preserves the cohesion of group A (though, not its position in the layout). The same occurs for group B (Fig. 3.14 (e)). However, the parameter change spreads B more than A – in image (e), we see that B maps to three groups, $B_1 \dots B_3$. These visualizations thus answer question (b). Group C behaves differently (Fig. 3.14 (f)): This group is split into two smaller groups C_1 and C_2 when we change our parameter. For question (c), thus, the answer is partially negative: not all groups are preserved in terms of spatial coherence upon parameter change.

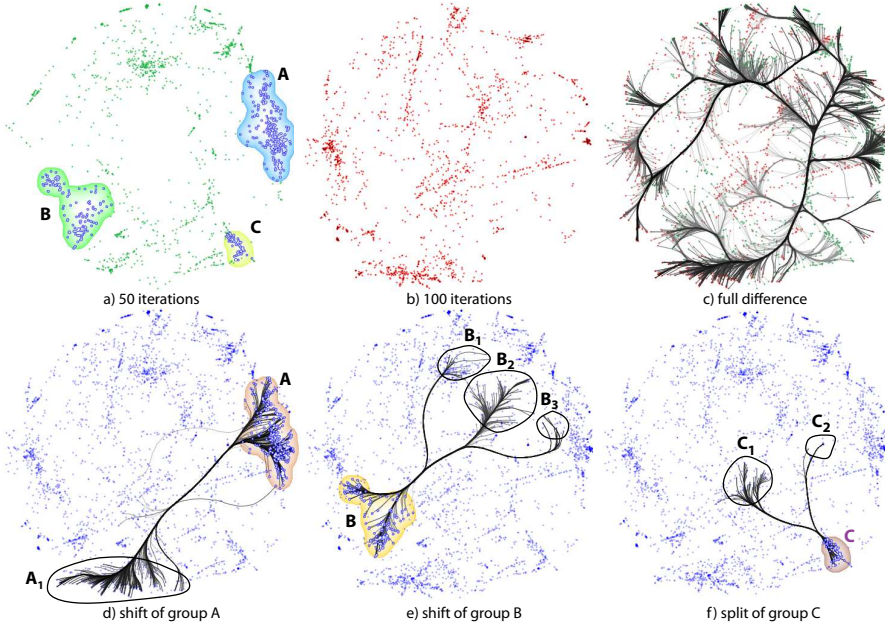


Figure 3.14: Applications – Shift between two LSP projections, for different numbers of force-directed iterations.

3.4 Discussion

We have implemented our visualization techniques in C++ using OpenGL 1.1, and tested them on Linux, Windows, and Mac OSX. Below we discuss several aspects of our techniques.

Computational scalability: For Delaunay triangulation and nearest-neighbor searches, we use the Triangle [167] and ANN [6] libraries. Both can handle over 100K points in subsecond time on a commodity PC. Further, we accelerate imaging operations using GPU techniques. For 2D distance transforms, we use the pixel-accurate Euclidean distance transform algorithm and GPU implementation proposed in [28]. On an Nvidia GT 330M, this allows us to compute shaded cushions and perform our Shepard interpolation at interactive frame rates for views of 1024^2 pixels. For edge bundling, we implemented KDEEB [85] fully on Nvidia’s CUDA platform. This yields a speed-up of over 30 times (on average) as compared to the C# implementation in [85] and allows bundling graphs of tens of thousands of edges in roughly one second. All in all, we achieve interactive querying and rendering of our views for projections up to 10K points.

Visual scalability: Our image-based approaches scale well to thousands of data points or more, even when little screen space is available. Moreover, all our techniques have a multiscale aspect: The parameters α and β (Eqns. 3.4, 3.5) effectively control the visual

scale at which we want to see false neighbors, missing neighbors, and the aggregate error. Increasing these values eliminates spatial outliers smaller than a given size, thereby emphasizing only coarse-scale patterns (see *e.g.* Fig. 3.1). The bundled views (Sec. 3.2.5) also naturally scales to large datasets given the inherent property of bundled edge layouts to emphasize coarse-scale connectivity patterns.

Genericity: Our visualizations are applicable to any DR algorithm, as long as one can compute an error distance matrix encoding how much 2D distances deviate from their n D counterparts (Eqn. 2.2). No internal knowledge of, or access to, the DR algorithms is needed – these can be employed as black boxes. This allows us to easily compare widely different DR algorithms, *e.g.* based on representatives, based on distance matrices, or based on direct use of the n D coordinates.

Ease of use: Our views are controlled by three parameters: α sets the scale of the visual outliers we want to show; β sets the radius around a point in which we want to display information, *i.e.*, controls the degree of space-filling of the resulting images; ϕ sets the percentage of most important missing neighbors we want to show. These parameters, as well as the interaction for selecting point groups (Sec. 3.2.6) are freely controllable by users by means of sliders and point-and-click operations.

Comparison: Similarly to Van der Maaten *et al.* [193], we use multiple views showing the same data points to explain a projection, *e.g.*, the false neighbors, missing neighbors view, missing neighbors finder, and group-related maps. However, the multiple maps in [193] are used to actually convey the projection, so the same point can have different locations and/or weights in different maps. In contrast, we use multiple views to convey different quality metrics atop of the same 2D projection, but keep the position of all observations the same in all these views. This simplifies the user’s task of correlating these multiple views, based on spatial positions. Similar to Aupetit [8], our error metrics encode discrepancies in distances in \mathbb{R}^n vs \mathbb{R}^2 . However, our error metrics are different. More importantly, our visualizations are different: Our false neighbors view does not show (a) spurious Voronoi cell edges far away from data points or (b) cell subdivision edges whose locations does not convey any information, since we (a) use distance-based blending and (b) continuous rather than constant per-cell interpolation (Sec. 3.2.3). Secondly, our missing neighbors finder (Sec. 3.2.5) can show one-to-many and many-to-many error relationships, whereas all other methods are constrained to one-to-one relationships. Finally, we can show errors at group level, whereas the other studied techniques confine themselves to showing errors at point level only.

Our projection comparison view (Sec. 3.2.7) is technically related to the method of Turkay *et al.*, which connects two 2D scatterplots to each other by lines linking their corresponding points [189]. However, Turkay *et al.* stress that line correspondences only work for a *small* number of points. In contrast, we use bundles to (a) show up to thousands of correspondences, and coloring and blending to encode correspondence importance.

Findings: It can be argued that our results are limited, as we did not decide, using our method, which of the studied DR algorithms are best. However, this was not the aim of our work. Rather, our goal was to present a set of visual techniques that help analyse the effect of parameters on projection quality for several DR techniques of interest. Deciding whether a certain degree of quality, *e.g.* in terms of false neighbors, missing neighbors, grouping problems, or projection stability is a highly context, dataset, and application-dependent task. Having such a context, our tools can be then used to assess (a) which are the quality problems, (b) how parameter settings affect them, and (c) whether these problems are acceptable for the task at hand. The same observation applies to the datasets used here. Our analyses involving these should be seen purely as test cases for assessing the quality problems of DR projections, and not as findings that affect the underlying problems captured by these datasets.

Distance vs neighborhood preservation: It can be argued that our error metrics, and corresponding visualizations, measure and respectively present a mix between the preservation of distances and neighborhoods: The aggregate error view is a ‘pure’ distance-related metric, very close to the well-known aggregated stress metric. The false neighbors view encodes distance errors into colors, but only for the Delaunay neighbors of each projected point. The missing neighbors view is also a pure distance-related metric. The missing neighbors finder highlights, indeed, the most important missing neighbors of a projected point, but uses a distance metric to define the notion of missing neighbors. The group analysis views lift the above interpretations at the level of a user-chosen group of projected points.

Obviously, distance preservation and neighborhood preservation are related notions: Perfect distance preservation implies perfect neighborhood preservation, as nearest neighbors are defined by distances. However, the converse is not true – neighborhoods are defined in terms of the order of points as given by their sorted set of distances, and not by the absolute value of distances. As such, very high distance-preservation errors are quite likely to cause also significant errors in neighborhood preservation, but the two types of errors are not reducible to each other. As all our error metrics proposed in this chapter are essentially distance-based, we refer to them as distance-preservation errors. In contrast, errors that use sets of nearest neighbors in their formulation, rather than distances, will be called neighborhood-preservation errors, and will be analysed separately in Chapter 4.

Limitations: A few limitations can be identified for the techniques presented in this chapter. The error metric described by Equation 3.1 can be quite sensitive to outliers. If there is an outlier in the original space (D^n) then $\max_{i,j} d^n(\mathbf{p}_i, \mathbf{p}_j)$ will be much larger than most other pairwise distances d^n . Hence, after naive linear normalization, the distances between all non-outlier points will be squeezed into a small portion of the output (normalized) space. Unless this outlier behavior is very well-represented in the projection (D^m), this will make the rightmost term of Equation 3.1 be quite different than the leftmost term even when they should be similar, so well-positioned points may still show errors. Considering also the phenomenon of ‘curse of dimensionality’, it can be very hard for projections to represent outliers well without crowding non-outlier points

in a very small area in the output space. The neighborhood-preservation error metrics proposed in the next chapter deal with this by considering, for each point i , only the ranks of neighbors of i , i.e., their discrete position in the sorted list of nearest-neighbors of i , and not their actual distances.

Regarding our visualizations, as outlined by the examples, they can show (a) which projection areas suffer from low quality; and (b) how two projections differ in terms of neighborhood preservation. However, we cannot directly explain (c) *why* a certain DR algorithm decided to place a certain point in some position; and (d) how the user should *tune* (if possible) the algorithm's parameters to avoid errors in a given area. In other words, we can explain the function $f : P$ (Eqn. 2.1) and its first derivatives over P , but not the inverse f^{-1} . This is a much more challenging task – currently not solved by any technique we know of. Further explaining such second-order effects to help users locally fine-tune a projection is subject to future work. Secondly, the parameter space P of some DR algorithms can be high-dimensional. So far, we can only analyse the variation of one or two parameters at a time. Extending this to several parameters is a second challenging next topic.

3.5 Conclusions

We have presented a set of visualization methods for the analysis of the quality of dimensionality-reduction (DR) algorithms in terms of their ability to preserve distances. We generically model such algorithms as functions from n D to 2D, parameterized in terms of the various settings of the respective projection algorithm. Next, we classify distance-related projection errors into false neighbors, missing neighbors, and aggregated projection error at both individual point and point-group level, and propose metrics to quantify these errors. We next propose several dense-pixel, visually scalable, techniques such as multi-scale scattered point interpolation and bundled edges to display out error metrics in ways that are visually and computationally scalable to large datasets and also work in a multiscale mode. We demonstrate our techniques by analysing the parameters of five state-of-the-art DR techniques.

In contrast to existing assessments of DR projections by aggregate figures, that can only infer overall precision, we offer more local tools to examine how neighborhoods and groups are mapped in the final projection. The usage of our techniques is simple and, most importantly, allows users of DR techniques to study their quality without needing to understand complex internal processes or the exact role of each parameter in the projections.

As noted in the introduction of this chapter, the techniques presented here only deal with distance errors in projections. While these are, naturally, important to quantify and see, they are not the only sources of problems in interpreting projections. As such, in Chapter 4, we will detail the issue of quantifying and visualizing neighborhood-preservation errors, the second most-important source of projection interpretation errors identified in Chapter 2. Separately, Chapter 5 will show how projection errors can be computed and used to assess distance preservation for 3D projections, and also compare the quality of 2D and 3D projections generated by the same projection technique.

After identifying the nature, distribution, and magnitude of errors in a given projection, a natural question for users is to assess which points get affected by such errors. This requires explaining groups of points in a projection in terms of the underlying high-dimensional variables. A related task of interest is to understand why errors appear for a given subset of points of a given dataset projected by a given algorithm. Both above tasks can be addressed by depicting the most important high-dimensional variables that cause a projection to place points close to, or far away from, each other. Techniques that address these tasks will be discussed in Chapter 6.

Contributions

The text of this chapter is based on the article “Visual analysis of dimensionality reduction quality for parameterized projections” (R. Martins, D. Coimbra, R. Minghim, A. Telea), *Computers & Graphics*, vol. 41, pp 26-42, 2014. The two first co-authors have had equal major contributions to this publications, and should be seen as joint first authors. Specific contributions of R. Martins involve: the proposal of dense image-based techniques to encode and visualize various forms of projection errors, based on smooth interpolation, and the design of the various error metrics (Sec. 3.2.2 and following); the adaptation and usage of edge bundles to visualize various types of errors (Secs. 3.2.5 and 3.2.6); and the selection, usage, and interpretation of projections constructed by the PLMP, Pekalska, and ISOMAP techniques (Sec. 3.3.2 and following).