Robust Spectral Clustering for Noisy Data

Modeling Sparse Corruptions Improves Latent Embeddings

Aleksandar Bojchevski Technical University of Munich a.bojchevski@in.tum.de Yves Matkovic Technical University of Munich matkovic@in.tum.de Stephan Günnemann Technical University of Munich guennemann@in.tum.de

ABSTRACT

Spectral clustering is one of the most prominent clustering approaches. However, it is highly sensitive to noisy input data. In this work, we propose a robust spectral clustering technique able to handle such scenarios. To achieve this goal, we propose a sparse and latent decomposition of the similarity graph used in spectral clustering. In our model, we jointly learn the spectral embedding as well as the corrupted data – thus, enhancing the clustering performance overall. We propose algorithmic solutions to all three established variants of spectral clustering, each showing linear complexity in the number of edges. Our experimental analysis confirms the significant potential of our approach for robust spectral clustering. Supplementary material is available at www.kdd.in.tum.de/RSC.

CCS CONCEPTS

•Computing methodologies →Machine learning approaches; Unsupervised learning; Spectral methods; •Information systems →Data mining; Clustering;

1 INTRODUCTION

Clustering is one of the fundamental data mining tasks. Among the variety of methods that have been introduced in the literature [1], spectral clustering [20] is one of the most prominent and successful approaches. It has been successfully applied in many domains ranging from computer vision to network analysis.

Since spectral clustering relies on a similarity graph only (e.g. connecting each instance with its m nearest neighbors), it is applicable to almost any data type, with vector data being the most frequent case. Spectral clustering embeds the data instances into a vector space that is spanned by the k eigenvectors corresponding to the k smallest eigenvalues of the graph's (normalized) Laplacian matrix. By clustering in this space, even complex structures can be detected – such as the half-moon data shown in Fig. 1 (left).

While spectral clustering is widely used in practice, one big issue is rarely addressed: it is highly sensitive to noisy input data. Fig. 1 illustrates this effect. While for the data on the left spectral clustering perfectly recovers the ground-truth clusters, the scenario on the right – with only slightly perturbed data – leads to a completely

KDD'17, August 13-17, 2017, Halifax, NS, Canada.



Figure 1: Spectral clustering (SC) is sensitive to noisy input. Left: SC detects the clustering. Right: SC fails. Our method (RSC) is successful in both scenarios.

wrong clustering for any of the three established versions [20] of spectral clustering. Spectral clustering fails in such scenarios.

In this work, we introduce a principle to robustify spectral clustering. The core idea is that the observed similarity graph is not perfect but corrupted by errors. Thus, instead of operating on the original graph – or performing some, often arbitrary, data cleaning that precedes the analysis – we assume the graph to be decomposed into two *latent* factors: the clean data and the corruptions. Following the idea that corruptions are sparse, we *jointly* learn the latent corruptions and the latent spectral embedding using the clean data.

For tasks such as regression [18], PCA [2], and autoregression [6, 8], such ideas have shown to significantly outperform non-robust techniques. And, indeed, also our method – called RSC – leads to clusterings that are more robust to corruptions. In Fig. 1 (right) our approach is able to detect the correct clustering structure. More precisely, our work is based on a sparse latent decomposition of the graph with the aim to optimize the eigenspace of the graph's Laplacian. This is in strong contrast to, e.g., robust PCA where the decomposition is guided by the eigenspace of the data itself. In particular, different Laplacians affect the eigenspace differently and require different solutions.

We note that the focus of this work is not on finding the number of clusters automatically. Principles using, e.g., the largest eigenvalue gap [14] might similarly be applied to our work. We left this aspect for future work. Overall, our contributions are:

- **Model:** We introduce a model for robust spectral clustering that handles noisy input data. Our principle is based on the idea of sparse latent decompositions. This is the first work exploiting this principle for spectral clustering, in particular tackling also the challenging case of normalized Laplacians.
- Algorithms: We provide algorithmic solutions for our model for all three established versions of spectral clustering using different Laplacian matrices. For our solutions we relate to principles such as Eigenvalue perturbation and the multidimensional Knapsack problem. In each case, the complexity of the overall method is linear in the number of edges.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

^{© 2017} ACM. ISBN 978-1-4503-4887-4/17/08...\$15.00.

DOI: http://dx.doi.org/10.1145/3097983.3098156

• **Experiments:** We conduct extensive experiments showing the high potential of our method, with up to 15 percentage points improvement in accuracy on real-world data compared to standard spectral clustering. Moreover, we propose two novel measures – *local purity* and *global separation* – which enable us to evaluate the intrinsic quality of an embedding without relying on a specific clustering technique.

2 PRELIMINARIES

We start with some basic definitions required in our work. Let A be a matrix, we denote with a_i the *i*-th *row*-vector of A and with $a_{i,j}$ the value at position i, j. A similarity graph is represented by a symmetric adjacency matrix $A \in (\mathbb{R}_{\geq 0})^{n \times n}$, with *n* being the number of instances. We denote the set of *undirected edges* as $\mathcal{E} = \{(i, j) \mid a_{i,j} > 0 \land i > j\}$. The set of edges incident to node *i* is given by $\mathcal{E}_i = \{(x, y) \in \mathcal{E} \mid x = i \lor y = i\}$. The vector representing the edges of A is written as $[a_{i,j}]_{(i,j) \in \mathcal{E}} = [a_e]_{e \in \mathcal{E}}$.

We denote with $d_i = \sum_j a_{i,j}$ the degree of node *i*, and with $D(A) = diag(d_1, \ldots, d_n)$ the diagonal matrix representing all degrees. We denote with *I* the identity matrix, whose dimensionality becomes clear from the context. Furthermore, as required for spectral clustering, we introduce different notions of Laplacian matrices: - unnormalized Laplacian: L(A) = D(A) - A

- normalized Laplacians: $L_{rw}(A) = D(A)^{-1}L(A)$

and $L_{sym}(A) = D(A)^{-1/2}L(A)D(A)^{-1/2}$

2.1 Spectral Clustering

Spectral clustering can be briefly summarized in three steps (see [20] for details). Step 1: Construct the similarity graph *A*. Different principles for the similarity graph construction exist. We focus on the symmetric *x*-nearest-neighbor graph, as it is recommended by [20] – any other construction can be used as well. Thus, the graph *A* is given by $a_{i,j} = 1$ if *i* is a *x* nearest neighbor of *j* or vice versa, and $a_{i,j} = 0$ else.

Step 2: Depending on the considered Laplacian, the next step is to compute the following eigenvectors¹:

- L(A): k first eigenvectors of L(A)

- $L_{rw}(A)$: k first generalized eigenv. of $L(A)u = \lambda D(A)u$

- $L_{sym}(A)$: k first eigenvectors of $L_{sym}(A)$

This step stems from the fact that spectral clustering tries to obtain solutions that minimize the ratio-cut/normalized-cut in the similarity graph. As shown in [20], an approximation to, e.g., the ratio-cut is obtained by the following trace minimization problem

$$\min_{\boldsymbol{H} \in \mathbb{R}^{n \times k}} Tr(\boldsymbol{H}^T \boldsymbol{L}(\boldsymbol{A})\boldsymbol{H}) \text{ subject to } \boldsymbol{H}^T \boldsymbol{H} = \boldsymbol{I}$$
(1)

The solution being the *k* first eigenvectors of the Laplacian *L* as stated above. Similar trace minimization problems can be formulated for the other Laplacians. We denote with $H \in \mathbb{R}^{n \times k}$ the matrix storing the eigenvectors as columns.

Step 3: Clustering on *H*. The spectral embedding of each instance *i* is given by the *i*-th row of *H*. To find the final clustering, the vectors h_i are (in case of L_{sym} first normalized and then) clustered using, e.g., *k*-means.

3 RELATED WORK

Multiple principles to improve spectral clustering have been introduced – focusing on different kinds of robustness. Surprisingly, many of the techniques [9, 11, 14, 23] are based on fully connected similarity graphs – even though nearest neighbor graphs are recommended [20]. First, using fully connected graphs highly increases the runtime – the considered matrices are no longer sparse – and, second, one has to select an appropriate scaling factor σ , required, e.g., for the Gaussian Kernel when constructing the graph (see [20]). Thus, many techniques [9, 11, 14, 23] focus on robustness regarding the parameter σ .

Local similarity scaling: [23] introduces a principle where the similarity is locally scaled per instance, i.e. the parameter σ changes per instance. By doing so, an improved similarity graph is obtained that better separates dense and sparse areas in the dataspace. The work [11] has extended this principle by using a weighted local scaling. The methods work well on noise-free data; however, they are still sensitive to noisy inputs.

Laplacian smoothing: [9] considers the problem of noisy data similar to our work, and they propose a principle of eigenvector smoothing. The initial Laplacian matrix is replaced by a smoothed version $M = \sum_{i=2}^{n} \frac{1}{\gamma + \lambda_i} \mathbf{x}_i \cdot \mathbf{x}_i^T$ where \mathbf{x}_i and λ_i are the eigenvectors/values of the original Laplacian matrix. Clustering is then performed on the eigenvectors of the matrix *M*. A significant drawback is that a full eigenvalue decomposition is required.

Data warping: [14] focuses on data where uniform noise has been added; not noisy data itself. They propose the principle of data warping. Intuitively, the data is transformed to a new space where noise points form its own cluster. Since they focus on fully connected graphs, noise can easily be detected by inspecting points with the lowest overall similarity. Since [9] and [14] are the most closely related works to our principle, we compare against them in our experiments.

Feature weighting: Focusing on a different scenario, multiple works have considered noisy/irrelevant features. In [10] a global feature weighting is learned in a semi-supervised fashion, thus, leading to an improved similarity matrix. [24] learns an affinity matrix based on random subspaces focusing on discriminative features. In [7], inspired by the idea of subspace clustering, feature weights are learned locally per cluster.

All the above techniques (except [7]) follow a *two-step*, *sequential* approach: They first construct an improved similarity graph/ Laplacian and then apply standard spectral clustering. In contrast, our method *jointly* learns the similarity graph and the spectral embedding. Both steps repeatedly benefit from each other.

Besides the above works focusing on general spectral clustering, different extended formulations have been introduced: [13] considers hypergraphs to improve robustness, [3] uses path-based characteristics. None of the techniques jointly learns a similarity matrix and the spectral embedding. Not focusing on robustness w.r.t. noise, [21] computes a doubly stochastic matrix by imposing low-rank constraints on the graph's Laplacian. It is restricted to the unnormalized Laplacian and leads to dense graphs, making it impractical for large data. Moreover, works such as [15] consider the problem of finding anomalous subgraphs using spectral principles, again not focusing on the case of noise.

 $^{^1 \}rm We$ denote with 'k first' eigenvectors, those k eigenvectors referring to the k smallest eigenvalues.

We further note that the spectral analysis is not restricted to a graph's Laplacian (as used in standard spectral clustering). The classical works of Davis-Kahan [5], for example, study the perturbation of a matrix X and the change of X's eigenspace. Following this line, [22] studies clustering based on the eigenspace of the adjacency matrix itself. In contrast, in this paper, we focus on the change of the eigenspace of L(X). In particular, we also consider the case of normalized Laplacians, which often lead to better results [20].²

4 ROBUST SPECTRAL CLUSTERING

In the following, we introduce the major principle of our technique – called RSC. For illustration purposes, we will start with spectral clustering based on the unnormalized Laplacian. The (more complex) principles for normalized Laplacians are described in Sec. 5

Let $A \in (\mathbb{R}_{\geq 0})^{n \times n}$ be the symmetric similarity graph extracted for the given data, with *n* being the number of instances in our data (see Sec. 2). Our major idea is that the similarity graph A is not perfect but might be corrupted (e.g. due to noisy input data). Any analysis performed on A might lead to misleading results.

Therefore, we assume that the observed graph A is obtained by two *latent* factors: A^c representing the corruptions and A^g representing the 'good' (clean) graph. More formally, we assume an additive decomposition³, i.e.

 $A = A^g + A^c$ with $A^g, A^c \in (\mathbb{R}_{>0})^{n \times n}$, both symmetric.

Instead of performing the spectral clustering on the corrupted A, our goal is to perform it on A^g . The core question is, how to find the matrices A^g and A^c ? In particular since clustering is an unsupervised learning task we don't know which entries in A might be wrong. For solving this challenge, we exploit two core ideas:

1) Corruptions are relatively rare – if they were not rare, i.e. the majority of the data is corrupted, a reasonable clustering structure can not be expected. Technically, we assume the matrix A^c to be sparse.

Let θ denote the maximal number of corruptions a user expects in the data. We require $||A^c||_0 \le 2 \cdot \theta$ where

$$||A^{c}||_{0} := |\{(i, j) | a_{i, j}^{c} \neq 0\}|$$

denotes the element-wise L_0 pseudo-norm ($2 \cdot \theta$ due to symmetry of the graph).

While θ constrains the number of corruptions *globally*, it is likewise beneficial to enforce sparsity *locally* per node. This can be realized by the constraint $\|a_i^g\|_0 \ge m$ for each node *i* (or equivalently: $\|a_i^c\|_0 \le |\mathcal{E}_i| - m$; we chose the first version due to easier interpretability: each node in A^g will be connected to at least *m* other nodes). Note that θ and *m* control different effects. To ignore either global or local sparsity, one can simply set the parameter to its extreme value ($\theta = \frac{1}{2} \|A\|_0$ or m = 1).

2) The detection of A^{g}/A^{c} is steered by the clustering process, i.e., we *jointly* perform the spectral clustering and the decomposition of A. This is in contrast to a sequential process where first the matrix is constructed and then the clustering is performed.



Figure 2: Spectral embeddings for data of Fig. 1 (right). Left: Spectral clustering; middle: RSC with $\theta = 10$, right: $\theta = 20$. RSC enhances the discrimination of points.

The strong advantage of a simultaneous detection is that we don't need to specify a separate – often arbitrary – objective for finding A^g , but the process is complete determined by the underlying spectral clustering. More precise, we exploit the equivalence of spectral clustering to trace minimization problems (see Sec. 2.1, Eq. (1)). Intuitively, the value of the trace in Eq. (1) corresponds to an approximation of the ratio-cut in the graph A. The smaller the value, the better the clustering. Thus, we aim to find the matrix A^g by minimizing the trace based on the Laplacian's eigenspace – subject to the sparsity constraints. Overall, our problem becomes:

PROBLEM 1. Given the matrix \mathbf{A} , the number of clusters k, the sparsity threshold θ , and the minimal number of nearest neighbors m per node. Find $\mathbf{H}^* \in \mathbb{R}^{n \times k}$ and $\mathbf{A}^{g^*} \in (\mathbb{R}_{\geq 0})^{n \times n}$ such that

$$(H^*, A^{g^*}) = \underset{H \ A^g}{\operatorname{argmin}} \operatorname{Tr}(H^T \cdot L(A^g) \cdot H)$$
(2)

subject to $\mathbf{H}^T \cdot \mathbf{H} = \mathbf{I}$ and $\mathbf{A}^g = \mathbf{A}^{gT}$ and $\|\mathbf{A} - \mathbf{A}^g\|_0 \le 2 \cdot \theta$ and $\|\mathbf{a}_i^g\|_0 \ge m \quad \forall i \in \{1, \dots, n\}$

The crucial difference between Eq. (1) and Problem 1 is that we now *jointly* optimize the spectral embedding H and the similarity graph A^g . The Laplacian matrix $L(A^g)$ is *no longer constant* but adaptive.

Figure 2 shows the strong advantage of this joint learning. Here, different spectral embeddings H (2nd and 3rd eigenvector since the 1st is constant) for the data in Fig. 1 (right) are shown. The left plot shows the embedding using *usual* spectral clustering. Due to the noisy input, the three groups are very close to each other and each spread out. Clustering on this embedding merges multiple groups and, thus, leads to low quality (for real-world data these embeddings look even harder as we will see in the experimental section). In contrast, the middle and right images show the spectral embedding learned by our technique when removing just 10 or 20 corrupted edges, respectively. Evidently, the learned embeddings highlight the clustering structure more clearly. Thus, by simultaneously learning the embedding and the corruptions, we improve the clustering quality.

4.1 Algorithmic Solution

While our general objective is hard to optimize (in particular due to the $\|.\|_0$ constraints the problem becomes NP-hard in general), we propose a highly efficient block coordinate-descent (alternating) optimization scheme to approximate it. That is, given H, we optimize for A^g/A^c ; and given A^g/A^c we optimize for H (cf. Algorithm 1). Of course, since A^c determines A^g and vice versa, it is sufficient to focus on the update of one, e.g., A^c . It is worth pointing out

²Surprisingly, many advanced spectral works still consider only the easier case of unnormalized Laplacians. Our competitors [9, 14] handle normalized Laplacians. ³This general decomposition not only leads to good performance, as we will see later, but also facilitates easy interpretation.

that in many works, the $\|.\|_0$ norm is simply handled by relaxation to the $\|.\|_1$ norm. In our work, in contrast, we aim to *preserve the interpretability* of the $\|.\|_0$ norm; for this, we derive a connection to the multidimensional Knapsack problem.

Update of *H*: Given A^c , the update of *H* is straightfarward. Since $A^g = A - A^c$ and therefore $L(A^g)$ are now constant, we can simply refer to Eq. (2): finding *H* is a standard trace minimization problem. The solution of *H* are the *k* first eigenvectors of $L(A^g)$.

Update of A^c : Clearly, since A^c needs to be non-negative, for all elements (i, j) with $a_{i,j} = 0$, it also holds $a_{i,j}^c = 0$. Thus, in the following, we only have to focus on the elements $a_{i,j}^c$ with $(i, j) \in \mathcal{E}$, i.e. the vector $[a_e^c]_{e \in \mathcal{E}}$. We base our update on the following lemma:

LEMMA 4.1. Given H, the solution for A^c minimizing Eq. (2) can be obtained by maximizing

$$f_1([a_e^c]_{e\in\mathcal{E}}) \coloneqq \sum_{(i,j)\in\mathcal{E}} a_{i,j}^c \cdot \left\| \boldsymbol{h}_i - \boldsymbol{h}_j \right\|_2^2$$
(3)

П

subject to the $\|.\|_0$ constraints and for each $e: a_e^c \in \{0, a_e\}$.

Proof. See appendix.

Exploiting Lemma 4.1, our problem can equivalently be treated as a set selection problem. For this, let $X \subseteq \mathcal{E}$ and $[v_e^X]_{e \in \mathcal{E}} = \boldsymbol{v}^X \in \mathbb{R}^{|\mathcal{E}|}$ be the vector with $v_e^X = \begin{cases} a_{i,j} & \text{if } (i,j) = e \in X \\ 0 & \text{else} \end{cases}$, our goal is

to find a set $X^* \subseteq \mathcal{E}$ maximizing $f_1(\boldsymbol{v}^{X^*})$ subject to the constraints. Accordingly, Problem 1 can be represented as (a special case of) a multidimensional Knapsack problem [16] operating on the set of edges \mathcal{E} :

COROLLARY 4.2. Given H. Let $X = \{e \in \mathcal{E} \mid x_e = 1\}$ be the solution of the following multidimensional Knapsack problem: Find $x_e \in \{0, 1\}, e \in \mathcal{E}$ such that $\sum_{e \in \mathcal{E}} x_e \cdot p_e$ is maximized subject to $\sum_{e \in \mathcal{E}} x_e \le \theta$ and $\forall i = 1, ..., n : \sum_{e \in \mathcal{E}_i} x_e \le |\mathcal{E}_i| - m$ where

$$p_e = p_{(i,j)} = a_{i,j} \cdot \left\| h_i - h_j \right\|_2^2$$
 (4)

The solution for A^c w.r.t. Eq. (2) corresponds to $\boldsymbol{v}^{\boldsymbol{\chi}}$.

This result matches the intuition of corrupted edges: The term p_e is high for instances whose embeddings are very dissimilar (i.e. they should not belong to the same cluster) but which are still connected by an edge.

While finding the optimal solution of a multidim. Knapsack problem is intractable, multiple efficient and effective approximate solutions exist [12, 16]. We exploit these approaches for our final algorithm. Following the principle of [12], we first sort the edges $e \in \mathcal{E}$ based on their ratio $p_e/\sqrt{s_e}$. Here, s_e is the number of constraints the variable x_e participates in. Since in our special case, *each* x_e participates in exactly three constraints, $s_e = 3$, it is sufficient to sort the edges based on the value p_e . We then construct a solution by adding one edge after another to A^c as long as the constraints are not violated. This approach leads to the best possible worst-case bound of $1/\sqrt{n+1}$ [12].

Algorithm 1 (lines 5-15) shows the update of A^c/A^g . Note that we do not need to sort the full edge set. It is sufficient to iteratively obtain the best edges. Thus, a priority queue PQ (e.g. a heap) is used (line 7, 10). The local $\|.\|_0$ constraints can simply be ensured by recording how many edges per node can still be removed (line

input :Similarity graph A, parameters k, θ , m**output**: Clustering C_1, \ldots, C_k 1 $A^g \leftarrow A;$ 2 while true do /* Update of H */ Compute Laplacian, matrix H, and trace; 3 if Trace could not be lowered then break; 4 /* Update of A^c/A^g */ $X = \emptyset$; 5 for each node *i* set $count_i \leftarrow |\mathcal{E}_i| - m$; 6 7 priority queue PQ on tuples (score, edge); for each edge $e \in \mathcal{E}$ add tuple (p_e, e) to PQ if $p_e > 0$ 8 [Eq. (4) or Eq. (6)]; while PQ not empty do 9 10 get first element from PQ \rightarrow (., $e_{best} = (i, j)$); if $count_i > 0 \land count_j > 0$ then 11 $X \leftarrow X \cup \{e_{best}\};$ 12 $count_i - -; count_j - -;$ 13

14 **if** $|X| = \theta$ **then** break;

15 construct A^c according to \boldsymbol{v}^{χ} ; $A^g = A - A^c$;

16 apply k-means on (normalized) vectors $(\mathbf{h}_i)_{i=1,...,n}$

Algorithm 1: Robust spectral clustering

6, 13). Thus, an edge can only be included in the result (line 12) if the incident nodes allow to do so (line 11).

The overall method for robust spectral clustering using unnormalized Laplacians iterates between the two update steps (lines 3-15). Note that in each iteration, line 8 considers all edges of the original graph. Thus, an edge marked as corrupted in a previous iteration might be evaluated as non-corrupted later. The algorithm terminates when the trace can not been improved further. In the last step (line 16), the *k*-means clustering on the improved *H* matrix is performed as usual.

Complexity: Using a heap, the update of A^c can be computed in time $O(|\mathcal{E}| + \theta' \cdot \log |\mathcal{E}|)$, where $\theta' \leq |\mathcal{E}|$ is the number of iterations of the inner while loop. Using power iteration, the eigenvectors H can be computed in time linear in the number of edges. Thus, overall, linear runtime can be achieved, as also verified empirically.

Sparse operations: It is worth mentioning that all operations performed in the algorithm operate on sparse data. This includes the computation of the Laplacian, its eigenvectors, and the constructions of A^c and A^g . Thus, even large datasets can easily be handled.

5 RSC: NORMALIZED LAPLACIANS

We now tackle the more complex cases of the two normalized Laplacians, which often lead to better clustering. For this, different algorithmic solutions are required.

5.1 Random Walk Laplacian

Spectral clustering based on L_{rw} corresponds to a *generalized* eigenvector problem using L [20]. Our problem definition becomes:

PROBLEM 2. Identical to Problem 1 but replacing the constraint $H^T \cdot H = I$ with $H^T \cdot D(A^g) \cdot H = I$.

Again, our goal is to solve this problem via block-coordinate descent. While the update of H is clear (corresponding to the first k generalized eigenvectors w.r.t. $L(A^g)$ and $D(A^g)$), using the

same approach for A^c/A^g as introduced in Sec. 4.1 turns out to be impractical: Since the constraint $H^T \cdot D(A^g) \cdot H = I$ now also depends on A^g , we get a highly restrictive constrained problem. As a solution, we propose a principle exploiting the idea of eigenvalue perturbation [19].

Using eigenvalue perturbation, we derive a matrix A^g aiming to minimize the sum of the k smallest generalized eigenvalues. Minimizing this sum is equivalent to minimizing the trace based on the normalized Laplacian's eigenspace. We obtain:

LEMMA 5.1. Given the eigenvector matrix **H** and the corresponding eigenvalues $\lambda = (\lambda_1, \dots, \lambda_k)$. An approximation of A^c minimizing the objective of Problem 2 can be obtained by maximizing

$$f_2([a_e^c]_{e\in\mathcal{E}}) = \sum_{(i,j)\in\mathcal{E}} a_{i,j}^c \left(\left\| \boldsymbol{h}_i - \boldsymbol{h}_j \right\|_2^2 - \left\| \sqrt{\boldsymbol{\lambda}} \circ \boldsymbol{h}_i \right\|_2^2 - \left\| \sqrt{\boldsymbol{\lambda}} \circ \boldsymbol{h}_j \right\|_2^2 \right)$$
(5)

subject to the $\|.\|_0$ constraints and for each $e: a_e^c \in \{0, a_e\}$.

Here, $\sqrt{.}$ denotes the element-wise square-root of the vector elements, and \circ the Hadamard product.

Clearly, the solution of the unnormalized case (Eq. (3)) and the normalized case (Eq. (5)) are structural very similar – and for solving it we can use the same principle as before (Algorithm 1), simply using as edge scores now the values

$$p_{e} = p_{(i,j)} = a_{i,j} \left(\left\| \boldsymbol{h}_{i} - \boldsymbol{h}_{j} \right\|_{2}^{2} - \left\| \sqrt{\lambda} \circ \boldsymbol{h}_{i} \right\|_{2}^{2} - \left\| \sqrt{\lambda} \circ \boldsymbol{h}_{j} \right\|_{2}^{2} \right)$$
(6)

Accordingly, also the complexity for finding A^c remains unchanged. Note that only edges with positive score need to be added to the queue (line 8).

Advantages: Comparing Eq. (6) with Eq. (4) one sees an additional 'penalty' term which takes the norm/length of the vectors h_i and h_j into account. Thereby, instances whose embeddings are far away from the origin get a lower (or even negative) score. This aspect is highly beneficial for spectral clustering: e.g., in the case of two clusters, the final clustering can be obtained by inspecting the sign of the 1d-embedding [20] – in general, intuitively speaking, clusters are separated by the origin (see Fig. 2 where the origin is in the center of the plots). Instances that are far away from the origin can be clearly assigned to their cluster; thus, marking their edges as corrupt might improve the clustering only slightly. In contrast, edges that are at the border between different clusters are the challenging ones – and exactly these are the ones preferred by Eq. (6).

5.2 Symmetric Laplacian

We now turn to the last case, spectral clustering using L_{sym} .

PROBLEM 3. Identical to Problem 1 but replacing Eq. (2) with

$$(H^*, A^{g^*}) = \underset{H, A^g}{\operatorname{argmin}} \operatorname{Tr}(H^T \cdot L_{sym}(A^g) \cdot H) \tag{7}$$

Using alternating optimization, the matrix H can easily be updated when A^g is given. For updating the matrix A^g (or equivalently A^c) we use the following result:

LEMMA 5.2. Given the eigenvector matrix \mathbf{H} . The matrix \mathbf{A}^c minimizing Eq. (7) can be obtained by maximizing

$$f_3([a_e^c]_{e \in \mathcal{E}}) \coloneqq \sum_{(i,j) \in \mathcal{E}} \frac{a_{i,j} - a_{i,j}^c}{\sqrt{d_i - d_i^c} \cdot \sqrt{d_j - d_j^c}} \cdot \mathbf{h}_i \cdot \mathbf{h}_j^T$$

subject to the $\|.\|_0$ constraints and $0 \le a_e^c \le a_e$, where $d_i^c = \sum_{e \in \mathcal{E}_i} a_e^c$.

PROOF. Similar to proof of Lemma 4.1; see appendix

What is the crucial difference between Lemma 5.2 and Lemma 4.1/5.1? For the previous solutions, the objective function has decomposed in independent terms. That is, when adding an edge to A^c , i.e. changing $a_{i,j}^c$ from 0 to $a_{i,j}$, the scores of the other edges are not affected. In Lemma 5.2, the sum in f_3 does *not* decompose into independent terms. In particular, the terms d_i^c in the denominator lead to a coupling of multiple edges.

While, in principle, f_3 can be optimized via projected gradient ascent, each gradient step would require to iterate through all edges. Therefore, as an alternative, we propose a more efficient greedy approximation: Similar to before, we focus on the solutions $\boldsymbol{v}^{\boldsymbol{X}}$. Starting with $\boldsymbol{X} = \boldsymbol{\emptyset}$, we iteratively let this set grow following a steepest ascent strategy. That is, we add the edge e_{best} to \boldsymbol{X} fulfilling

$$e_{best} = \arg\max_{e \in \mathcal{E}'} f_3(\boldsymbol{v}^{\mathcal{X} \cup \{e\}}) \tag{8}$$

where \mathcal{E}' indicates the edges that could be added to \mathcal{X} without violating the constraints. Naively computing Equation (8) requires $|\mathcal{E}'| \cdot |\mathcal{E}|$ many steps – and since we perform multiple iterations to let \mathcal{X} grow, it results in a runtime complexity of $O(\theta \cdot |\mathcal{E}|^2)$; obviously not practical. In the following, we show how to compute this result more efficiently.

Definition 5.3. Let $\mathcal{X} \subseteq \mathcal{E}$, $d_i^{\mathcal{X}} := d_i - \sum_{e \in \mathcal{E}_i \cap \mathcal{X}} a_e$, and $p_{i,j} := a_{i,j} \cdot \mathbf{h}_i \cdot \mathbf{h}_j^T$. We define $s(i, w, \mathcal{X}) := \sum_{i=1}^{n} \sum_{e \in \mathcal{E}_i \cap \mathcal{X}} \left(\frac{1}{1 - 1} - \frac{1}{1 - 1} \right) p_{i,j}$

$$s(i, w, X) := \sum_{\substack{j \\ (i,j) \in \mathcal{E}_i \setminus X \\ \forall (j,i) \in \mathcal{E}_i \setminus X}} \left(\frac{1}{\sqrt{d_i^X - w} \sqrt{d_j^X}} - \frac{1}{\sqrt{d_i^X} \sqrt{d_j^X}} \right)^{p_{i,j}}$$

for each node *i*, and

$$\delta(e, X) := \left(\frac{1}{\sqrt{d_i^X}\sqrt{d_j^X}} - \frac{1}{\sqrt{d_i^X - a_e}\sqrt{d_j^X}} - \frac{1}{\sqrt{d_i^X}\sqrt{d_j^X - a_e}}\right)p_{i,j}$$

for each edge e = (i, j), and

$$\Delta(e, X) := s(i, a_e, X) + s(j, a_e, X) + \delta(e, X)$$

COROLLARY 5.4. Given X and $\mathcal{E}' \subseteq \mathcal{E} \setminus X$. It holds

$$\arg\max_{e\in\mathcal{E}'}f_3(\boldsymbol{v}^{\mathcal{X}\cup\{e\}}) = \arg\max_{e\in\mathcal{E}'}\Delta(e,\mathcal{X})$$

PROOF. See appendix.

By exploiting Corollary 5.4, we can find the best edge according to Eq. (8), by only considering the terms $\Delta(e, X)$. This term can be interpreted as the gain in f_3 when adding the edge e to the set X. After computing the scores s(i, w, X) for each *node*, $\Delta(e, X)$ can be evaluated in constant time per *edge*.

Moreover, let e = (i, j), for each non-incident edge $(i', j') = e' \in \mathcal{E} \setminus (\mathcal{E}_i \cup \mathcal{E}_j)$ it obviously holds $s(i', w, X) = s(i', w, X \cup \{e\})$ and $\delta(e', X) = \delta(e', X \cup \{e\})$. Thus, assume the edge $e_{best} = (i, j)$ has been identified and added to X. For finding the *next* best edge, only the scores s(i, ., .) and s(j, ., .) need to be updated; followed by an evaluation of δ for all edges incident to the nodes i and j. The remaining nodes and edges are not affected; their s, δ , and Δ values are unchanged.



Figure 3: Spectral embedding of banknote data based on L_{sym} . Note that the dataset contains *two* clusters. Left: Standard spectral clustering; middle & right: Our method ($\theta = 10$ and 20). The learned embeddings increase the discrimination between the points. The two clusters stand out more clearly.

Exploiting these results, we compute the set X similar to Algorithm 1 (lines 5 - 15): Initially, compute for each node *i* and unique edge weight $a_{i,j}$ the term $s(i, a_{i,j}, X)$. Then compute for each edge *e* the term $(\Delta(e, X), e)$ and add it to the priority queue PQ. These steps can be done in time $O(\gamma \cdot |\mathcal{E}|)$, where γ is the number of unique edge weights per node. Within the while loop: Every time the best element $e_{best} = (i, j)$ from the PQ is retrieved, we recompute s(i, ., X) and s(j, ., X), followed by a recomputation of $\delta(e, X)$ for all incident edges. Noticing that there are at most $2 \cdot x$ many incident edges (*x* nearest-neighbor graph) these steps can be done in time $O(\gamma \cdot x + x \cdot log(|\mathcal{E}|))$.

Overall, this leads to a time complexity of $O(\gamma \cdot |\mathcal{E}| + \theta \cdot (x \cdot log(|\mathcal{E}|) + \gamma \cdot x))$. Note that the worst case (each edge has a unique weight) corresponds to $\gamma = x$. In this case we obtain $O(x \cdot |\mathcal{E}| + \theta \cdot (x \cdot log(|\mathcal{E}|) + x^2))$. For our case of spectral clustering using nearest-neighbor graphs, however, it holds $\gamma = 1$. In this case, we obtain an algorithm with complexity

$O(|\mathcal{E}| + \theta \cdot x \cdot log(|\mathcal{E}|))$

Thus, being linear in the number of edges.

In summary, the principle for solving Eq. (7) is almost identical to Algorithm 1 with the additional overhead of re-evaluating the term $\Delta(e, X)$ for the edges incident to e_{best} . The full pseudocode of this algorithm and the detailed complexity analysis are provided in the supplementary material for convenience.

6 EXPERIMENTS

Setup. We compare our method, called RSC, against spectral clustering (SC), and the two related works AHK [9] and NRSC [14]. We denote with RSC-L_{xy} the different variants of our method using the corresponding Laplacian. For all techniques, we set the number of clusters *k* equal to the number of clusters in the data. As default values we construct nearest neighbor graphs with 15 neighbors, allowing half of the edges to be removed per node ($m = 0.5 \cdot x$). While [14] uses a principle for automatically setting their parameters, the obtained results were often extremely low. Thus, we manually optimized their parameters to obtain better solutions. All experiments are averaged over several k-means runs to ensure stability. All used datasets are publicly available/on our website. Real world data: We use handwritten digits (pendigits; 7494 instances; 16 attributes; 10 clusters)⁴, banknote authentication data (1372 inst.; 5 att.; 2 clus.)⁴, iris (150 inst.; 4 att.; 3 clus.)⁴, and USPS data (9298 inst.; 256 att.;

10 clus.)⁵. Further, we use two random subsamples of the MNIST data (10k/20k inst., 784 att., 10 clus.) because our competitors can not handle larger samples due to their cubic complexity. Synthetic data: Besides the well known moon data as shown in Fig. 1, where the vectors' positions are perturbed based on Gaussian noise using different variance, we also generate synthetic similarity graphs based on the planted partitions model [4]: Given the clusters, we randomly connect each node to *x* percent of the other nodes in its cluster. Additionally, we add a certain fraction of noise edges to the graph. By default we generate data with 1000 instances, x = 0.3 and 20 clusters. We evaluate the clustering quality of the different approaches using NMI (1=best). We start with an in-depth analysis of our technique followed by a comparison with competing techniques.

Spectral embedding. RSC optimizes the spectral embedding H by learning the matrix A^g . Thus, we start by analyzing the spectral embeddings obtained by RSC. In Fig. 2 we illustrated the spectral embeddings for the data of Fig. 1 (right). Standard spectral clustering fails on this data, since the embedding (left plot in Fig. 2) leads to unclear groupings. In contrast, applying our technique, we obtain the embeddings as shown in Fig. 2 (right): the three clusters stand out; thus, perfect clustering structure can be obtained.

A similar behavior can be observed for real world data. Fig. 3 shows the spectral embedding of the banknote data (*two* clusters) regarding L_{sym} (the other Laplacians show similar results). On the left we see the original embedding: The points do not show a clear separation in two groups. In the middle and right plot, we applied RSC with θ =10 and θ =20, respectively. As shown, the separation between the points clearly increases. The embedding gets optimized leading to higher clustering accuracy. As we will see later, for the banknote data, the NMI score increases from 0.46 to 0.61.

Sparsity threshold. As indicated in Fig. 3, increasing the sparsity threshold might lead to a clearer separation. We now analyze this aspect in more detail. Figure 4 (left) analyzes a two-moons datasets with noise of 0.1. We vary θ for all three techniques. $\theta = 0$ corresponds to original spectral clustering using the corresponding Laplacian; clearly, its quality is low. As shown, for all techniques we observe an increase in the clustering quality until a stable point is reached. Fig. 4 (right) shows the same behavior for the banknote data. The removal of corrupted edges improves the clustering results. All three variants are able to reach the highest NMI of 0.61.

⁴https://archive.ics.uci.edu/ml/

⁵http://www.cs.nyu.edu/~roweis/data.html



clustering quality. Left: two moons data; right: banknote.

Remark: Using the variant based on *L*, at some point, the quality will surely drop again. When all corrupted edges have been removed, one will start to remove 'good' edges. The reason is that the terms in Eq. (3) are always non-negative. In contrast, using L_{rw}/L_{sum} (Eq. (3); Cor. 5.4, Δ), edges connecting points within the same cluster will often obtain negative scores. Those edges will *never* be included in the matrix A^c – independent of θ . Thus, in general, the later two versions are more robust regarding θ .

According to our definitions we aim to minimize the trace. In Fig. 5 we illustrate the value of the trace for the setting of Fig. 4 (right). Since the trace between the different Laplacian can not be meaningfully compared in absolute values, we plot it relative to the trace obtained by standard spectral clustering. For all of our approaches the trace can successfully be lowered by a significant amount; thus, confirming the effectiveness of our learning approach. Note also that our algorithms often need only around 10 iterations to converge to these good results.

Detection of corrupted edges. Next, we analyze how well our principles are able to spot corrupted edges. For this, we artificially added corrupted edges to the similarity graph based on the planted partition model. We used two different settings: in one case 10% of all edges in the graph are corrupted; in the other even 20% of all edges. Knowing the corrupted edges, we measure the precision $p = |B \cap A|/|B|$ and recall $r = |B \cap A|/|A|$, where A denotes the corrupted edges, and B the edges removed by our technique.

Fig. 6 shows the results when increasing the number of removed edges (i.e. θ). For the 10% noise case (left plot), we observe a very high precision which stays at the optimal value until 1200 - only the corrupted edges are removed. Note that the absolute number of corrupted edges in the data is 1261. Likewise, the recall is continuously increasing until around 0.96. Thus only a few corrupted edges could not be detected. The scenario with 20% noise (3605

1

0,8

Precision (all variants)

corrupted edges) is more challenging. While L_{sym} obtains a result very close to optimal, L_{rw} and L perform slightly worse. Thus, for these techniques also some 'good' edges get removed. Note that the curves do not need to be monotonic. Due to the joint optimization, different edges can be removed for each parameter setting.

Overall, for realistic scenarios of noise, all techniques perform well – with L_{sum} often being the best one.

Robustness. In the next experiment, we analyze the robustness of our method regarding perturbed data. That is, we study how an increasing degree of noisy data effects the clustering quality. We refer to the established moon data and perturb it randomly according to Gaussian noise with variance increased from 0 to 0.115. To highlight the variation in the clustering quality we average the results over 10 datasets for each noise parameter.

Fig. 7 shows the results of our principles and standard spectral clustering. The lines represents the mean NMI, while the error bars represent the variance. Note that for standard SC we report the best result among all three Laplacian (for each dataset individually). Thus, standard spectral clustering gets an additional strong benefit. Clearly, spectral clustering is not robust and rapidly decreases in quality. Interestingly, for the moon data, L performs best. In any case, all of our approaches clearly outperform the baseline.

Comparison of Runtime. We now turn our attention to the comparison between RSC and related techniques. First, we briefly evaluate the runtime behaviour. The experiments were conducted on 2.9 GHz Intel Core i5 with 8GB of RAM running Matlab R2015a. Fig. 8 shows the overall runtime for each method on pendigits. To obtain larger data, we performed supersampling; adding small noise (variance of 0.1) to avoid duplicates. Confirming our complexity analysis, RSC scales linear in the number of edges - and it easily handles graphs with around 1 mio edges. Not surprisingly, standard spectral clustering is the fastest. The competing techniques are much slower due to their cubic complexity in the number of nodes; they can only handle small graphs. For the larger datasets, they did not finish within 24 hours.

Comparison of clustering quality. Next, we provide an overview of the clustering quality. For all techniques we used the symmetric normalized Laplacian since it performed best. Even though our main aim is to improve spectral clustering approaches, we additionally report the results of two famous clustering principles: (A) k-means and (B) density-based clustering (here: mean shift). For the later, we tuned the bandwidth parameter to obtain highest scores. As already mentioned in the set-up, the competing techniques' parameters were tuned as well. For RSC, we simply used



Figure 5: Our method obtains better (lower) trace values (trace of SC=100%)



0,8

Figure 6: Our method achieves high precision and recall; the corrupted edges are successfully detected. Left: 10% noise; right: 20% noise.

0.7



Figure 7: Robustness to noise. Our RSC clearly outperforms spectral clustering.



Figure 8: Runtime analysis. RSC scales linear in the number of edges

RSC L



Figure 9: Robustness of all techniques on banknote. RSC is very stable.

a very large θ and let the method automatically decide how many edges to mark as corrupted – one advantage of L_{sym} (see remark in experiment on sparsity threshold). Table 1 summarizes the results for the different datasets.

Besides using	data	SC	NRSC	АНК	(A)	(B)	RSC
the full datasets.	moone	0.47	0.00	0.53	0.10	0.34	1 00
wa usa tha prin	banknoto	0.46	0.77	0.55	0.17	0.03	0.61
we use the prin-	LISDS	0.40	0.47	0.52	0.05	0.05	0.01
ciple of [9, 14]	MNIST-10 K	0.70	0.85	0.77	0.01	0.15	0.03
and additionally	MNIST 20 K	0.71	0.76	0.70	0.40	0.40 d n f	0.75
select sets of the	iris	0.70	0.70	0.71	0.40	0.72	0.70
data More pre-	nendigits	0.82	0.83	0.55	0.70	0.72	0.82
· 1 C 1	nendigits-16	0.86	0.87	0.82	0.88	0.00	0.02
cisely, from the	pendigits-146	0.00	0.07	0.00	0.88	0.01	0.91
pendigits data we	T T		0.71			1	0.50
select specific dig-	Table 1: Clustering quality						

its indicated with pendigits-xyz.

As shown, in many cases our technique outperforms the competing techniques. In some scenarios by even 15 percentage points w.r.t. spectral clustering. Though, it is also fair to mention that not for all datasets an improvement can be achieved. Our method clearly outperforms k-means and density based clustering and it finishes for all these datasets in a few seconds to minutes. In contrast, NRSC and AHK required already around one and three hours respectively on the larger MNIST data.

Comparison of robustness. Next, we analyze the robustness of the methods by artificially adding noise to the real data. To ensure that the cluster detection is indeed getting more difficult, we specifically add corruptions to the similarity graph connecting different clusters. The results for the banknote data are presented in Fig. 9. As shown, at the beginning all techniques remain at their quality level obtained on the original data, with RSC obtaining the highest quality. Adding more corruptions, however, standard spectral clustering drops very quickly and sharply to low quality. In contrast, RSC stays at its highest level for the longest time. AHK is quite stable as well, while NRSC is much more sensitive.

Comparison of the embeddings' quality. One of our main hypothesis is that jointly learning the embedding and the corruptions leads to improved embeddings. Thus, lastly, we study the quality of the embeddings learned by all techniques. While we have already seen in Table 1, that the NMI scores of RSC are good, such measures give only limited insights how the underlying embedding space looks like. How can we measure the quality of the embeddings? In particular, we aim to derive statistics that do not depend on applying a clustering technique on the data - instead we want to

evaluate the embedding based on the ground truth classes⁶ only. We argue that two properties should be fulfilled: (a) Local purity. In a good embedding, the instances within a local neighborhood should belong to the same class. (b) Global separation: In a good embedding, it should be possible to distinguish between instances of different classes by inspecting the intra-class and inter-class distances only. That is, the classes should be easily separable.

Evaluation of Local Purity. Let h_i denote the embedding of instance *i* and $c_i \in C$ its class according to the ground truth. We define the purity $pur_x(i)$ of the neighborhood around instance i (x-nearest neighbors) as the largest fraction of instances belonging to the same class. Formally: Let $NN_x(i)$ denote the set of x nearest neighbors of *i* in the embedding space and $occ_x(c, i)$ the number of times the class *c* occurs in the neighborhood of node *i* (including the node itself), i.e. $occ_x(c, i) = |\{j \in NN_x(i) \cup \{i\} \mid c_i = c\}|$. Then the purity is given by

$$pur_x(i) = \frac{1}{x+1} \max_{c \in C} occ_x(c, i)$$

The overall local purity (at scale x) is defined as the average over all instances:

$$PUR(x) = \frac{1}{N} \sum_{i=1}^{N} pur_{x}(i)$$

In the best case PUR(x) = 1, each neighborhood contains instances of a single class only; in the worst case 1/k where k = |C| is the number of classes. As an example, imagine an embedding that looks like Figure 1: We would observe good purity (PUR(x)=1) for small x. Slightly increasing x, the purity decreases since different classes get merged, demonstrating the low quality of such an embedding.

Results: Fig. 10 shows the result for banknote and USPS for all competing techniques. *x* is scaled from 1 to the maximal cluster size to ensure that all scales are captured. As a baseline we also evaluate the purity of the original input data (i.e. the embedding space is the raw data). As seen, in both plots, the original data only has good purity for small x, but drops quickly. This indicates complex shapes like in Fig. 1. For banknote, RSC has consistently the highest local purity. The embedding well reflects the ground truth locally. SC and AHK perform slightly worse. NRSC, in contrast, drops very quickly similar to the baseline. For USPS, all techniques clearly outperform the baseline. While for small x RSC is slightly below

⁶We specifically use the term 'class' to indicate the groups given by the ground truth, not the groups detected by an arbitrary clustering method applied on the embedding.



Figure 10: Evaluation of local purity (left: banknote, right: USPS). RSC's embedding represents the classes well locally.

the competitors, it outperforms them for larger x, which is also reflected by a higher NMI.

Evaluation of Global Separation. We formalize global separation by extending the idea of the Silhouette coefficient [17]. For each class c we compute the list of pairwise distances $P_{c,c}$ of all instances within the class, as well as the list of pairwise distances $P_{c,c'}$ between instances from class c to c', i.e.

$$P_{c,c'} = [dist(\boldsymbol{h}_i, \boldsymbol{h}_j)]_{i \in C_c, j \in C_{c'}}$$

where $C_c = \{i \mid c_i = c\}$ is the set of all instances from class *c*.

For each list we compute the average over the $(x \cdot 100)\%$ smallest elements, denoted as $P_{c,c'}(x)$. Following the Silhouette coefficient, we then compute the difference between the *within class* distances and the distance to the *closest other class*, i.e.

$$GS_{c}(x) = \frac{P_{c,c'}(x) - P_{c,c}(x)}{\max\{P_{c,c'}(x), P_{c,c}(x)\}}$$

where $c' = \arg \min_{c' \neq c} P_{c,c'}(x)$. In the best case $GS_c(x) = 1$, in the worst case -1. $GS_c(x)$ can intuitively be regarded as a robust extension of the Silhouette coefficient (w.r.t. the ground truth classes). For x = 1, it resembles the Silhouette coefficient w.r.t. class c. For x < 1 only parts of the distances are considered, thus, capturing that the embedding might not completely represent the ground truth. Imagine an embedding that resembles Fig. 1, $GS_c(x)$ will be relatively low due to similar inter-class and intra-class distances.

Results: Fig. 11 shows the result for two exemplary classes. An overview of all classes is available in the supp. material. Again, the raw data shows the worst result, with scores consistently below 0.25 indicating no good separation/clusteredness of the class labels in the space. In contrast, RSC obtains extremely high scores in both datasets up to a very high x: for banknote until 0.85, for USPS 0.97. That is, a very large fraction of the ground truth class is well separated and clustered in the learned embedding. Clearly, not every class from the data shows such perfect result since the NMI scores are 0.61 and 0.85. The (sharp) drops at the end indicate that some of the instances of the ground truth class are wrongly assigned to a different region in the embedding space. When trying to include these (x = 1), the score highly drops. The competing approaches consistently perform worse, showing no good match between the ground truth and the clusteredness of the embedding.

The results on the USPS data also indicate that *local purity* and *global separation* of an embedding are indeed two different properties: While SC and AHK have good results on the local purity, they perform poor regarding global separation. The learned embeddings of RSC capture well both properties confirming the benefit of our joint learning principle.



Figure 11: Evaluation of global separation (left: banknote, right: USPS). RSC's embedding separates the classes well.

7 CONCLUSION

We proposed a spectral clustering technique for noisy data. Our core idea was to decompose the similarity graph into two latent factors: sparse corruptions and clean data. We jointly learned the spectral embedding as well as the corrupted data. We proposed three different algorithmic solutions using different Laplacians. Our experiments have shown that the learned embeddings clearly emphasize the clustering structure and that our method outperforms spectral clustering and state-of-the-art competitors.

Acknowledgments. This research was supported by the German Research Foundation, Emmy Noether grant GU 1409/2-1, and by the Technical University of Munich - Institute for Advanced Study, funded by the German Excellence Initiative and the European Union Seventh Framework Programme under grant agreement no 291763, co-funded by the European Union.

REFERENCES

- C. C. Aggarwal and C. K. Reddy. Data clustering: algorithms and applications. CRC Press, 2013.
- [2] E. J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? JACM, 58(3):11, 2011.
- [3] H. Chang and D.-Y. Yeung. Robust path-based spectral clustering. *Pattern Recognition*, 41(1):191–203, 2008.
- [4] A. Condon and R. M. Karp. Algorithms for graph partitioning on the planted partition model. *Random Structures and Algorithms*, 18(2):116–140, 2001.
- [5] C. Davis and W. M. Kahan. The rotation of eigenvectors by a perturbation. iii. SIAM SINUM, 7(1):1-46, 1970.
- [6] N. Günnemann, S. Günnemann, and C. Faloutsos. Robust multivariate autoregression for anomaly detection in dynamic product ratings. In WWW, pages 361–372, 2014.
- [7] S. Günnemann, I. Farber, S. Raubach, and T. Seidl. Spectral subspace clustering for graphs with feature vectors. In *ICDM*, pages 231–240, 2013.
- [8] S. Günnemann, N. Günnemann, and C. Faloutsos. Detecting anomalies in dynamic rating data. In KDD, pages 841–850, 2014.
- [9] H. Huang, S. Yoo, H. Qin, and D. Yu. A robust clustering algorithm based on aggregated heat kernel mapping. In *ICDM*, pages 270–279, 2011.
- [10] F. Jordan and F. Bach. Learning spectral clustering. Adv. Neural Inf. Process. Syst, 16:305-312, 2004.
- [11] W. Kong, S. Hu, J. Zhang, and G. Dai. Robust and smart spectral clustering from normalized cut. *Neural Computing and Applications*, 23(5):1503–1512, 2013.
- [12] D. Lehmann, L. I. Oćallaghan, and Y. Shoham. Truth revelation in approximately efficient combinatorial auctions. *Journal of the ACM*, 49(5):577–602, 2002.
- [13] X. Li, W. Hu, C. Shen, A. Dick, and Z. Zhang. Context-aware hypergraph construction for robust spectral clustering. *TKDE*, 26(10):2588–2597, 2014.
- [14] Z. Li, J. Liu, S. Chen, and X. Tang. Noise robust spectral clustering. In ICCV, pages 1–8, 2007.
- [15] B. A. Miller, M. S. Beard, P. J. Wolfe, and N. T. Bliss. A spectral framework for anomalous subgraph detection. *Trans. on Signal Proc.*, 63(16):4191–4206, 2015.
- [16] J. Pfeiffer and F. Rothlauf. Analysis of greedy heuristics and weight-coded eas for multidimensional knapsack problems and multi-unit combinatorial auctions. In *GECCO*, pages 1529–1529, 2007.
- [17] P. J. Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Computational and applied mathematics*, 20:53–65, 1987.
- [18] P. J. Rousseeuw and A. M. Leroy. Robust regression and outlier detection, volume 589. John Wiley, 2005.

- [19] G. Stewart and J.-G. Sun. Matrix Perturbation Theory. Academic Press Boston, 1990.
- [20] U. von Luxburg. A tutorial on spectral clustering. Statistics and Computing, 17(4):395-416, 2007.
- [21] X. Wang, F. Nie, and H. Huang. Structured doubly stochastic matrix for graph based clustering. In KDD, pages 1245–1254, 2016.
- [22] L. Wu, X. Ying, X. Wu, and Z. Zhou. Line orthogonality in adjacency eigenspace with application to community partition. In *IJCAI*, pages 2349–2354, 2011.
- [23] L. Zelnik-Manor and P. Perona. Self-tuning spectral clustering. In NIPS, pages 1601–1608, 2004.
- [24] X. Zhu, C. Loy, and S. Gong. Constructing robust affinity graphs for spectral clustering. In CVPR, pages 1450–1457, 2014.

APPENDIX

PROOF OF LEMMA 4.1. Note that $L(A^g) = D(A^g) - A^g = (D(A) - D(A^c)) - (A - A^c) = L(A) - L(A^c)$. Thus, Eq. (2) can equivalently be written as $\operatorname{Tr}(H^T \cdot L(A^g) \cdot H) = \operatorname{Tr}(H^T \cdot (L(A) - L(A^c)) \cdot H) = \operatorname{Tr}(H^T \cdot L(A) \cdot H) - \operatorname{Tr}(H^T \cdot L(A^c) \cdot H)$.

Given H, the term $\operatorname{Tr}(H^T \cdot L(A) \cdot H)$ is constant. Thus, minimizing the previous term is equivalent to maximizing $\operatorname{Tr}(H^T \cdot L(A^c) \cdot H)$.

Let \boldsymbol{y}_k be a column vector of \boldsymbol{H} . Noticing that (see [20]) $\boldsymbol{y}_k^T \boldsymbol{L}(\boldsymbol{A}^c) \boldsymbol{y}_k$ = $\sum_{i,j} \frac{1}{2} \cdot a_{i,j}^c \cdot (y_{k,i} - y_{k,j})^2$, and exploiting the orthogonality of \boldsymbol{H} it follows: Tr $(\boldsymbol{H}^T \cdot \boldsymbol{L}(\boldsymbol{A}^c) \cdot \boldsymbol{H}) = \sum_k \sum_{i,j} \frac{1}{2} \cdot a_{i,j}^c \cdot (y_{k,i} - y_{k,j})^2 = \sum_{i,j} \frac{1}{2} \cdot a_{i,j}^c \cdot \|\boldsymbol{h}_i - \boldsymbol{h}_j\|_2^2$, where the last step used $y_{k,i} = h_{i,k}$.

To ensure that A^c as well as A^g are non-negative, it holds $0 \le a_{i,j}^c \le a_{i,j}$. Thus, if $a_{i,j} = 0$ then $a_{i,j}^c = 0$. Exploiting this fact and the symmetry of the graph leads to $\sum_{i,j} \frac{1}{2} \cdot a_{i,j}^c \cdot \|\mathbf{h}_i - \mathbf{h}_j\|_2^2 = \sum_{i,j \in \mathcal{E}} a_{i,j}^c \cdot \|\mathbf{h}_i - \mathbf{h}_j\|_2^2$.

$$\begin{split} &\sum_{(i,j)\in\mathcal{E}} a_{i,j}^c \cdot \left\| \boldsymbol{h}_i - \boldsymbol{h}_j \right\|_2^2. \\ &\text{Next, we show that there exists a solution where each } a_{i,j}^c \in \{0, a_{i,j}\}. \\ &\text{As known, } 0 \leq a_{i,j}^c \leq a_{i,j}. \text{ Let } M = [a_e^c]_{e\in\mathcal{E}} \text{ be a maximum of Eq.} \\ &(3) \text{ where some } a_{i,j}^c > 0 \text{ but } < a_{i,j}. \text{ Let } M' \text{ be the solution where } \\ &\text{this entry is replaced by } a_{i,j}^c = a_{i,j}. \text{ Since only } \|.\|_0 \text{ constraints are } \\ &\text{used, } M \text{ and } M' \text{ fulfill the same constraints. Since } \left\| \boldsymbol{h}_i - \boldsymbol{h}_j \right\|_2^2 \text{ is non-negative, } f_1(M') \geq f_1(M). \text{ It follows, that a solution minimizing Eq. (2) \\ &\text{can be found by investigating } a_{i,j}^c = 0 \text{ or } a_{i,j}^c = a_{i,j} \text{ only.} \\ &\square \end{split}$$

PROOF OF LEMMA 5.1. The goal is to find a matrix A^g whose sum of the first *k* eigenvalues is minimal (and fulfills the given constraints). Since, however, A^g is not known, we refer to the principle of eigenvalue perturbation.

Let A^t be the matrix obtained in the previous iteration of the alternating optimization and let y_i be the *i*-th generalized eigenvector of $L(A^t)$ (these are the *columns* of the matrix H from above, i.e. $y_{i,j} = h_{j,i}$). Furthermore, denote the corresponding eigenvalues with λ_i . We define $L(A^g) - L(A^t) = \Delta L$ and $D(A^g) - D(A^t) = \Delta D$.

Based on the theory of eigenvalue perturbation [19], the eigenvalue λ_i^g of $L(A^g)$ can be approximated by

$$\lambda_i^g \approx \lambda_i + \boldsymbol{y}_i^T \cdot (\Delta L - \lambda_i \cdot \Delta D) \cdot \boldsymbol{y}_i$$

= $\lambda_i + \boldsymbol{y}_i^T \cdot ((L(A^g) - L(A^t)) - \lambda_i \cdot (D(A^g) - D(A^t))) \cdot \boldsymbol{y}_i$
ing the fact that $L(A^g) = L(A) - L(A^g)$ and $D(A^g) = D(A^g)$

Using the fact that $L(A^g) = L(A) - L(A^c)$ and $D(A^g) = D(A) - D(A^c)$, and after rearranging the terms, we obtain

$$\lambda_i^g \approx \overbrace{\lambda_i + \boldsymbol{y}_i^T \cdot ((\boldsymbol{L}(\boldsymbol{A}) - \boldsymbol{L}(\boldsymbol{A}^t)) - \lambda_i \cdot (\boldsymbol{D}(\boldsymbol{A}) - \boldsymbol{D}(\boldsymbol{A}^t))) \cdot \boldsymbol{y}_i}^{-\boldsymbol{U}_i} - \underbrace{\boldsymbol{y}_i^T \cdot ((\boldsymbol{L}(\boldsymbol{A}^c)) - \lambda_i \cdot (\boldsymbol{D}(\boldsymbol{A}^c))) \cdot \boldsymbol{y}_i}_{=:g_i}$$

Since c_i is constant, minimizing λ_i^g is equivalent to maximizing g_i . Simplifying yields:

$$g_i = \boldsymbol{y}_i^T \cdot \boldsymbol{L}(\boldsymbol{A}^c) \cdot \boldsymbol{y}_i - \lambda_i \cdot \boldsymbol{y}_i^T \cdot \boldsymbol{D}(\boldsymbol{A}^c) \cdot \boldsymbol{y}_i$$

$$= \sum_{j,j'} \frac{1}{2} a_{j,j'}^{c} (y_{i,j} - y_{i,j'})^{2} - \lambda_{i} \sum_{j} y_{i,j}^{2} \cdot d_{j}^{c}$$

where $d_{j}^{c} = [D(A^{c})]_{j,j} = \sum_{j'} a_{j,j'}^{c}$. Thus
 $g_{i} = \sum_{j,j'} \frac{1}{2} a_{j,j'}^{c} (y_{i,j} - y_{i,j'})^{2} - \lambda_{i} y_{i,j}^{2} a_{j,j'}^{c}$
 $= \sum_{j,j'} a_{j,j'}^{c} \left(\frac{1}{2} (y_{i,j} - y_{i,j'})^{2} - \lambda_{i} y_{i,j}^{2}\right)$

and exploiting the symmetry of the graph, we obtain

$$g_{i} = \sum_{(j,j')\in\mathcal{E}} a_{j,j'}^{c} \left((y_{i,j} - y_{i,j'})^{2} - \lambda_{i} y_{i,j}^{2} - \lambda_{i} y_{i,j'}^{2} \right)$$

Since the overall goal is to minimize $\sum_{i=1}^{k} \lambda_i^g$, we aim at maximizing

$$\sum_{i=1}^{k} g_{i} = \sum_{i=1}^{k} \sum_{(j,j')\in\mathcal{E}} a_{j,j'}^{c} \left((y_{i,j} - y_{i,j'})^{2} - \lambda_{i} y_{i,j}^{2} - \lambda_{i} y_{i,j'}^{2} \right)$$
$$= \sum_{(j,j')\in\mathcal{E}} a_{j,j'}^{c} \left(\sum_{i=1}^{k} (y_{i,j} - y_{i,j'})^{2} - \sum_{i=1}^{k} \lambda_{i} y_{i,j}^{2} - \sum_{i=1}^{k} \lambda_{i} y_{i,j'}^{2} \right)$$

By noticing that $y_{i,j} = h_{j,i}$ we obtain

$$=\sum_{(j,j')\in\mathcal{E}}a_{j,j'}^{c}\left(\underbrace{\left\|\boldsymbol{h}_{j}-\boldsymbol{h}_{j'}\right\|_{2}^{2}-\left\|\sqrt{\lambda}\circ\boldsymbol{h}_{j}\right\|_{2}^{2}-\left\|\sqrt{\lambda}\circ\boldsymbol{h}_{j'}\right\|_{2}^{2}}_{x}\right)$$

Note that some of the terms *x* might be negative. Clearly, since we aim to maximize the equation – and since $a_{i,j}^c \ge 0$ – for these terms we have to choose $a_{i,j}^c = 0$. For the remaining (non-negative) terms, the same arguments apply as in the proof of Lemma 4.1: i.e. they are either 0 or $a_{i,j}$. Thus, overall, for each term we have $a_e^c \in \{0, a_e\}$.

PROOF OF LEMMA 5.2. Note that $a_{i,j}^g = a_{i,j} - a_{i,j}^c$ and $d_i^g = d_i - d_i^c$. Let \mathbf{y}_k be a column vector of \mathbf{H} . It holds $\mathbf{y}_k^T \cdot \mathbf{L}_{sym}(\mathbf{A}^g)\mathbf{y}_k \stackrel{[20]}{=} \sum_{i,j} \frac{1}{2} a_{i,j}^g \left(\frac{y_{k,i}}{\sqrt{d_i^g}} - \frac{y_{k,j}}{\sqrt{d_j^g}}\right)^2 = \sum_{i,j} \frac{1}{2} a_{i,j}^g \left(\frac{y_{k,i}^2}{d_i^g} + \frac{y_{k,j}^2}{\sqrt{d_j^g}} - \frac{2\cdot y_{k,i}y_{k,j}}{\sqrt{d_i^g}\sqrt{d_j^g}}\right) = \sum_i \frac{1}{2} y_{k,i}^2 + \sum_j \frac{1}{2} y_{k,j}^2 - \sum_{i,j} \frac{a_{i,j}^g y_{k,i}y_{k,j}}{\sqrt{d_i^g}\sqrt{d_j^g}}$. Since \mathbf{y}_k is given, the first two terms are constant. Furthermore, due to orthogonality it holds $Tr(\mathbf{H}^T \mathbf{L}_{sym}\mathbf{H}) = \sum_k \mathbf{y}_k^T \cdot \mathbf{L}_{sym}\mathbf{y}_k$. Thus, minimizing the trace is equivalent to maximizing

$$\sum_{k} \sum_{i,j} \frac{a_{i,j}^{g} y_{k,i} y_{k,j}}{\sqrt{d_{i}^{g}} \sqrt{d_{j}^{g}}} = \sum_{i,j} \frac{a_{i,j}^{g}}{\sqrt{d_{i}^{g}} \sqrt{d_{j}^{g}}} \mathbf{h}_{i} \cdot \mathbf{h}_{j}^{T}, \text{ noticing that } y_{k,i} = h_{i,j}.$$
Evaluating the graph's symmetry concludes the proof

Exploiting the graph's symmetry concludes the proof.

PROOF OF COROLLARY 5.4. Adding e = (i, j) to X has the following effects: the term a_e^c changes from 0 to a_e ; the degree of the two incident nodes becomes $d_i^{X \cup \{e\}} = d_i^X - a_e$. Therefore,

$$f_{3}(\boldsymbol{v}^{X\cup\{e\}}) = f_{3}(\boldsymbol{v}^{X}) - \frac{p_{e}}{\sqrt{d_{i}^{X}} \cdot \sqrt{d_{j}^{X}}} - \sum_{\substack{(x, y) \in (\mathcal{E}_{i} \cup \mathcal{E}_{j}) \setminus X \\ (x, y) \neq (i, j)}} \frac{p_{x, y}}{\sqrt{d_{x}^{X}} \cdot \sqrt{d_{x}^{X}}} + \sum_{\substack{x \neq j \\ (i, x) \in \mathcal{E}_{i} \setminus X \\ \forall (x, i) \in \mathcal{E}_{i} \setminus X \\ \forall (x, i) \in \mathcal{E}_{i} \setminus X}} \frac{p_{i, x}}{\sqrt{d_{i}^{X} - a_{e}} \sqrt{d_{x}^{X}}} + \sum_{\substack{x \neq j \\ (j, x) \in \mathcal{E}_{j} \setminus X \\ \forall (x, j) \in \mathcal{E}_{j} \setminus X}} \frac{p_{x, j}}{\sqrt{d_{j}^{X} - a_{e}} \sqrt{d_{x}^{X}}} = f(\boldsymbol{r}^{X}) + a(i, q, Y) + b(q, X)$$

 $= f_3(\boldsymbol{v}^{\Lambda}) + s(i, a_e, X) + s(j, a_e, X) + \delta(e, X) = f_3(\boldsymbol{v}^{\Lambda}) + \Delta(e, X)$

Since X is given, $f_3(\boldsymbol{v}^X)$ is constant. Thus, the edge $e \in \mathcal{E}'$ maximizing $f_3(\boldsymbol{v}^{X \cup \{e\}})$ is found by maximizing $\Delta(e, X)$.