

Spectral Image Segmentation with Global Appearance Modeling

Jeova F. S. Rocha Neto* Pedro F. Felzenszwalb†

October 7, 2022

Abstract

We introduce a new spectral method for image segmentation that incorporates long range relationships for global appearance modeling. The approach combines two different graphs, one is a sparse graph that captures spatial relationships between nearby pixels and another is a dense graph that captures pairwise similarity between *all pairs of pixels*. We extend the spectral method for Normalized Cuts to this setting by combining the transition matrices of Markov chains associated with each graph. We also derive an efficient method for sparsifying the dense graph of appearance relationships. This leads to a practical algorithm for segmenting high-resolution images. The resulting method can segment challenging images without any filtering or pre-processing.

1 Introduction

Image segmentation is a fundamental problem in computer vision. Spectral clustering methods pioneered by the normalized cuts approach [1] provide simple and powerful algorithms based on fundamental graph-theoretic notions and computational linear algebra.

Spectral clustering methods are formulated using an objective function defined by a graph. The classical constructions used for image segmentation focus on pairwise similarity between nearby pixels. In this paper we introduce a new spectral method that incorporates long range relationships for global appearance modeling. The resulting method can segment challenging images without any filtering or pre-processing. Figure 1 shows several results obtained with the proposed method. Figure 6 shows how the new method significantly outperforms the original normalized cuts formulation.

We use a dense graph to capture the global appearance of regions. We show the normalized cut criteria in this graph reflects the distributions of pixel values in each region using a kernel density estimate. The measure penalizes the overlap between distributions in different regions.

To implement our image segmentation approach we extend the normalized cuts spectral algorithm to a setting where there are multiple graphs that encode different grouping cues. Our approach for image segmentation combines two graphs. We provide a natural interpretation for the normalized cut criteria on each of these graphs. One of the graphs is sparse and does not depend on the image data, it simply captures spatial relationships between pixels. The other graph is dense and captures pairwise similarity between *all pairs of pixels*, irrespective of their spatial location.

The direct implementation of spectral methods to segment high resolution images is challenging due to high memory and computational requirements. We tackle this challenge using a graph sparsification approach that enables the efficient segmentation of high resolution images.

We show experimental results with a variety of images and provide a quantitative evaluation using a dataset of synthetic images with Brodatz textures. Our approach achieves highly accurate results in this setting despite the complex appearance of the textures.

*Department of Computer Science, Haverford College, Haverford, PA, USA

†School of Engineering, Brown University, Providence, RI, USA



Figure 1: Segmentation results using the proposed method

2 Previous work

The mathematical foundation of spectral graph partitioning can be dated back to [2] and [3]. The seminal work of Shi and Malik [1] built on this foundation to develop the normalized cuts method for image segmentation. Since then a significant body of work has emerged in the fields of both spectral clustering and image segmentation (see, i.e. [4]). An important development on the theoretical understanding of normalized cuts can be found in [5], where the authors traced the parallel between the normalized cuts criteria and low conductivity sets in Markov chains.

On the application and implementation side, the normalized cuts algorithm poses two main challenges when applied to image segmentation: how to (efficiently) construct a graph over the set of pixels and how to perform the spectral decomposition. In order to solve them, the usual approach is to construct sparse graphs where each pixel is connected to nearby pixels [1, 6, 7]. Furthermore, works such as [8, 9, 10] reduce the computation burden of spectral segmentation by relying on approximate eigenvector solvers. Our method instead starts with a dense graph connecting all pairs of pixels and uses an edge sparsification technique based on importance sampling. We also provide an importance sampling implementation of the original similarity graph proposed in [1].

A key step in developing spectral segmentation algorithms is to decide which information to extract from the image in order to compute pixel affinities (the edge weights). While the initial implementations used raw intensities or filter-banks [1, 5, 11], further developments adopted the intervening contour cue from probability of boundary (Pb) maps [12, 6, 7]. In this paper, we show that using raw pixel intensities (or color) is enough to obtain satisfactory segmentation results in complex images. In particular, our method is able to outperform traditional filter-bank based methods in segmentation of textured images.

Our spectral segmentation algorithm is also related to multiview clustering methods. In multiview clustering one is concerned on clustering the data using different features (or views). One of the earlier results on this problem can be traced to [13], where the authors propose multiview counterparts of traditional clustering algorithms, such as EM and K -means.

In the spectral multiview realm, the work presented in [14] proposes a two-view clustering algorithm that computes the normalized eigenvectors arising from a bipartite graph that encodes the views. The authors in [15] describe an iterative procedure that enforces consistency among views by solving the generalized eigenvalue problem for each view separately using the eigenvectors from other views computed in previous iterations. Finally, other relevant approaches to multiview clustering involve combining graphs and weights arising from different views into a sole graph either via a convex combination [16, 17] or according to their Laplacian’s power mean [18]. To the best of our knowledge, no multiview clustering algorithm has been developed for image segmentation. One possible explanation for this is the computational burden of existing multi-view clustering methods.

Considering image segmentation broadly speaking, there is an increasing popularity of deep learning based solutions, some of which find their inspiration on normalized cuts [19, 20]. We refer the interested reader to the work in [21] for a comprehensive survey on these techniques. Departing from this learning based paradigm, our method does not require training data and can be applied in different settings with little or no fine-tuning. Moreover, our method produces interpretable results. Furthermore, the proposed algorithm demonstrates the power and practical use of long range relationships between pixels, which are not, or at least not explicitly, present within typical deep learning frameworks.

3 Background

3.1 Graph cuts and spectral clustering

Let $G = (V, E, w)$ be an undirected weighted graph. A cut (A, B) is a partition of V into two disjoint sets. We will consider the weight of a cut in different graphs that have the same set of vertices. Let $w(i, j) = 0$ when $\{i, j\} \notin E$. The weight of a cut (A, B) in G is defined as,

$$\text{Cut}(A, B|G) = \sum_{i \in A, j \in B} w(i, j).$$

In the context of clustering and image segmentation it is typical to use large weights to indicate that elements are similar and should not be separated. In this case we can look for the *minimum cut* to find an optimal partition of V . However, this strategy is heavily biased towards imbalanced cuts, such as having a single node on one side. This motivated the introduction of the celebrated normalized cut criteria and algorithm [1].

The *normalized cut* value is defined as,

$$\text{NCut}(A, B|G) = \frac{\text{Cut}(A, B|G)}{\text{Vol}(A|G)} + \frac{\text{Cut}(A, B|G)}{\text{Vol}(B|G)} = \text{Vol}(V|G) \frac{\text{Cut}(A, B|G)}{\text{Vol}(A|G) \text{Vol}(B|G)}.$$

Here $\text{Vol}(A|G)$ is a measure of the ‘‘volume’’ of A defined as $\text{Vol}(A|G) = \sum_{i \in A, j \in V} w(i, j)$.

The spectral algorithm introduced in [1] solves a continuous relaxation of the minimum NCut problem. Let W be the weighted adjacency matrix of G . Let D be the diagonal degree matrix with $D(i, i) = \sum_{j \in V} W(i, j)$. The matrix $L = D - W$ is the *Laplacian* of G .

The NCut algorithm solves a generalized eigenvector problem,

$$Lx = \lambda Dx. \tag{1}$$

The algorithm selects the eigenvector x with second smallest eigenvalue, and partitions V by thresholding x . In [5] the NCut criteria and algorithm is described in terms of a Markov chain. Let $P = D^{-1}W$. The matrix P is the transition matrix of a Markov chain over the vertices V . The long term behavior of this Markov chain can be characterized by the solutions of the eigenvector problem,

$$Px = \lambda x. \tag{2}$$

A solution (λ, x) to the eigenvector problem in (2) leads to a solution $(1 - \lambda, x)$ to the generalized eigenvector problem in (1) and vice-versa. Therefore the generalized eigenvector x used in the NCut algorithm corresponds to the eigenvector of P with second largest eigenvalue.

3.2 Traditional Graph Construction for Image Segmentation

The classical application of normalized cuts for image segmentation involves a graph H where the vertices represent the image pixels and the weights reflect simultaneously the appearance similarity and distance between pairs pixels.

Let $H = (V, E, w)$ be a graph where the vertices V are the pixels in an image and the edges E connect every pair of pixels. We use $I(j)$ to denote the appearance (such as the brightness or color) of pixel j . We use $X(j)$ to denote the spatial location of the same pixel. Now define,

$$w(i, j) = \exp\left(-\frac{\|I(i) - I(j)\|^2}{2\sigma_I^2}\right) \exp\left(-\frac{\|X(i) - X(j)\|^2}{2\sigma_X^2}\right). \tag{3}$$

The graph H combines two grouping cues in a single real valued weight¹ and shares similarities with bilateral filtering [22]. Using the normalized cut criteria, pixels are encouraged to be grouped together if

¹The graph defined here differs slightly from the one used in [1] because in [1] the weight of an edge is set to 0 if the distance between i and j is above a threshold.

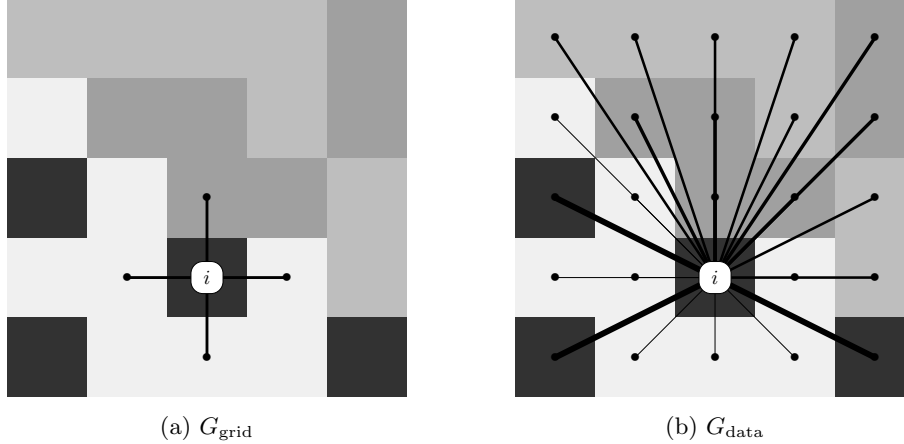


Figure 2: The edges connecting to a pixel i in G_{grid} and G_{data} . In the each image, the thickness of each link represents the weight of the edge connecting a pair of pixels.

they have similar appearance *and* are close to each other. Note, however, that pixels that have similar appearance but are far away are not encouraged to be grouped because the corresponding weight is close to zero. Similarly, neighboring pixels that have very different appearance, such as in a textured region, are also not encouraged to be grouped.

4 New Criteria for Image Segmentation

We combine two normalized cut values to obtain a new criteria for image segmentation. We break the grouping cues (spatial proximity and appearance similarity) into two separate graphs, G_{grid} and G_{data} . Both graphs are defined over the same set of vertices, corresponding to the pixels in an image.

1. The graph G_{grid} is a grid over the image pixels, where each pixel is connected to the four neighboring pixels with an edge of weight 1. This graph encourages neighboring pixels to be grouped together, independent of their appearance.
2. The graph G_{data} is a fully connected graph that encourages pixels with similar appearance to be grouped together, independent of their location. The weights in G_{data} are based on appearance similarity of pixels, and do not depend on pixel locations,

$$w(i, j) = \exp\left(-\frac{\|I(i) - I(j)\|^2}{2\sigma^2}\right).$$

Figure 2 illustrates the graph construction of both G_{grid} and G_{data} on a image with two regions.

4.1 Spatial Information: G_{grid}

Let (A, B) be a cut in the grid graph. The cut defines a segmentation of the image into two regions, with a boundary Γ between them. The cut value, $\text{Cut}(A, B|G_{\text{grid}})$, counts the number of neighboring pixels that are in different regions. In general the cut value in the grid graph and similar graphs can be seen as a measure of the length of the boundary Γ (see [23]).

Observation 1.

$$\text{Cut}(A, B|G_{\text{grid}}) \approx \text{Len}(\Gamma).$$

This is a commonly used measure of spatial coherence in image segmentation problems (see, e.g., [24]). Although the criteria $\text{Cut}(A, B|G_{\text{grid}})$ leads to spatially coherent segmentations and is widely used in practice, it gives most preference to trivial solutions with a small (single pixel) region.

Using the previous observation and noting that $\text{Vol}(S|G_{\text{grid}}) \approx 4|S|$ for $S \subseteq V$ we can derive an expression for the value of a normalized cut in the grid graph.

Observation 2.

$$\text{NCut}(A, B|G_{\text{grid}}) \approx \frac{|V| \text{Len}(\Gamma)}{4 |A||B|}.$$

Minimizing this criteria encourages solutions where the boundary Γ between the two regions is short (to minimize $\text{Len}(\Gamma)$) and where the two regions have similar size (to maximize $|A||B|$).

4.2 Global Appearance Information: G_{data}

Now we consider the weight of cuts and normalized cuts in G_{data} .

For $S \subseteq V$ we use g_S to denote a kernel density estimate defined by the values of pixels in S ,

$$g_S(c) = \frac{1}{|S|} \sum_{i \in S} K(I(i) - c).$$

Proposition 1.

$$\text{Cut}(A, B|G_{\text{data}}) = (2\pi\sigma^2)^{\frac{d}{2}} |A||B| \langle g_A, g_B \rangle.$$

Here d is the dimension of $I(j)$ ($d = 1$ for graylevel images and $d = 3$ for RGB images), g_A and g_B are density estimates defined using a Gaussian kernel, and $\langle \cdot, \cdot \rangle$ denotes the standard inner product of functions.

Proof. We use the fact that the convolution of two Gaussians with equal variance is a Gaussian with twice the variance,

$$\begin{aligned} \sum_{i \in A, j \in B} w_{\text{data}}(i, j) &= \sum_{i \in A} \sum_{j \in B} \exp\left(-\frac{\|I(i) - I(j)\|^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{\frac{d}{2}} \sum_{i \in A} \sum_{j \in B} \int_{-\infty}^{\infty} \left(\frac{1}{\pi\sigma^2}\right)^d \exp\left(-\frac{\|I(i) - c\|^2}{\sigma^2}\right) \exp\left(-\frac{\|I(j) - c\|^2}{\sigma^2}\right) dc \\ &= (2\pi\sigma^2)^{\frac{d}{2}} \int_{-\infty}^{\infty} \left(\sum_{i \in A} K(I(i) - c)\right) \left(\sum_{j \in B} K(I(j) - c)\right) dc \\ &= (2\pi\sigma^2)^{\frac{d}{2}} |A||B| \int_{-\infty}^{\infty} g_A(c) g_B(c) dc = (2\pi\sigma^2)^{\frac{d}{2}} |A||B| \langle g_A, g_B \rangle. \end{aligned}$$

□

The proposition above is related to the Laplacian PDF Distance in [25]. It is also related to the work in [26] where a different graph construction was used to define global appearance models.

The weight of a cut in G_{data} will be minimized when the pixel values in the two regions have complementary support. Although this intuitively makes sense, the measure encourages regions to be unbalanced in size due to the term $|A||B|$ multiplying $\langle g_A, g_B \rangle$.

In order to derive an expression for $\text{NCut}(A, B|G_{\text{data}})$, we first use a similar reasoning as in the proposition above to note that $\text{Vol}(S|G_{\text{data}}) = (2\pi\sigma^2)^{(d/2)} |S||V| \langle g_S, g_V \rangle$. Then, from the definition of the normalized cut we obtain the following result.

Proposition 2.

$$\text{NCut}(A, B|G_{\text{data}}) = \langle g_V, g_V \rangle \frac{\langle g_A, g_B \rangle}{\langle g_A, g_V \rangle \langle g_B, g_V \rangle}.$$

This criteria is minimized when the distributions g_A and g_B have little overlap and both have significant overlap with g_V . In particular it penalizes solutions where one region does not represent a significant amount of the image data.

4.3 Combining Spatial and Appearance information

The normalized cut values in G_{grid} and G_{data} provide complementary measures for image segmentation. To combine the spatial and appearance cues we use a convex combination,

$$\text{MixNCut}(A, B) = (1 - \lambda) \text{NCut}(A, B|G_{\text{data}}) + \lambda \text{NCut}(A, B|G_{\text{grid}}).$$

The parameter $\lambda \in [0, 1]$ controls the relative importance of the two normalized cut measures.

We interpret $\text{MixNCut}(A, B)$ as a mixture of an appearance and a spatial term,

$$\text{MixNCut}(A, B) \approx (1 - \lambda) \left(\langle g_V, g_V \rangle \frac{\langle g_A, g_B \rangle}{\langle g_A, g_V \rangle \langle g_B, g_V \rangle} \right) + \lambda \left(\frac{|V| \text{Len}(\Gamma)}{4 |A| |B|} \right).$$

The first term encourages a partition of the image into regions with dissimilar color distributions, while the second term encourages a spatially coherent partition. Both terms are normalized and avoid biases towards solutions with small regions. Note that each term is normalized in a particular way that is natural and has appropriate dimensions for the individual measures.

5 Segmentation Algorithm

Let G_1 and G_2 be two weighted graphs over the same set of vertices. Now we describe a spectral method we have used as a heuristic for optimizing a convex combination of two normalized cut values,

$$\text{MixNCut}(A, B|G_1, G_2) = (1 - \lambda) \text{NCut}(A, B|G_1) + \lambda \text{NCut}(A, B|G_2).$$

Our approach is based on the Markov chain and conductance interpretation of normalized cuts described in [5]. Let W_1 and W_2 be the weighted adjacency matrices of the two graphs while D_1 and D_2 are the diagonal degree matrices. Let,

$$P_1 = D_1^{-1}W_1,$$

$$P_2 = D_2^{-1}W_2$$

$$P = (1 - \lambda)P_1 + \lambda P_2. \tag{4}$$

The matrices P_1 and P_2 define two Markov chains on V . The matrix P also defines a Markov chain on V where in one step we follow P_1 with probability $(1 - \lambda)$ and P_2 with probability λ . We compute the second largest eigenvector of P to find a cut (A, B) with small conductance.

In our experiments, we use a Lanczos Process to compute the second largest eigenvector of P . We use K -means with $K = 2$ to cluster the entries in the eigenvector into 2 clusters.

5.1 Graph Sparsification

When the matrix P is sparse we can compute the required eigenvector much more quickly. The grid G_{grid} is sparse but G_{data} is dense. We sparsify the graph using a random sampling approach.

The approach described here is complementary to other methods that have been used to speed up the computation of eigenvectors for clustering. One such method is based on Nystrom approximation [9]. Another approach involves power iteration [10].

Let G be a weighted graph. To construct a sparse graph G' we independently sample m edges (with replacement) from G , with probabilities proportional to the edge weights. The weight of each sampled edge

is set to 1 (adding up weights if there is repetition). With this approach the expected value of a cut (A, B) in G' equals the value of the cut in G up to a scaling factor of $(m/\text{Vol}(V|G))$. Moreover, if m is sufficiently large then with high probability every cut in G' has weight close to the cut value in G (up to a scaling factor of $(m/\text{Vol}(V|G))$) (see, e.g., [27]). For the experiments in this paper we use $m = \alpha|V|$ with $\alpha > 0$ to obtain a sparse graph approximating G .

To implement this approach efficiently for G_{data} we need to sample edges with probability proportional to their weights $w(i, j)$ *without* enumerating all possible edges. We use an importance sampling method as a practical alternative.

First we partition V into L (≈ 1000 in practice) sets S_1, \dots, S_L with low appearance variance. We do this greedily, starting with a single set and repeatedly partitioning the set with highest variance into two using the K -means algorithm. Let m_i be the mean appearance of pixels in S_i and

$$q(a, b) = |S_a||S_b| \exp\left(-\frac{\|m_a - m_b\|^2}{2\sigma^2}\right).$$

To sample an edge for G' we select a random pair S_a and S_b with probability proportional to $q(a, b)$. We then select $i \in S_a$ and $j \in S_b$ uniformly at random. Finally, we add the edge $\{i, j\}$ to G' with weight $w'(i, j) = |S_a||S_b|w(i, j)/q(a, b)$.

6 Numerical Experiments

The algorithms in this section were implemented in MATLAB and run on a computer with an Intel i5-6200U CPU @ 2.30GHz using 8 GB of RAM running Linux.

6.1 Segmentation accuracy measure

To measure the accuracy of a segmentation we use the Jaccard Index to compare pairs of regions [28],

$$J(S, Q) = |S \cap Q|/|S \cup Q|.$$

Let I be an image with a ground-truth binary segmentation (S, Q) . Let (A, B) be the result of a binary segmentation algorithm. To evaluate the accuracy of (A, B) we consider (1) pairing S with A and Q with B or (2) pairing S with B and Q with A . We define the accuracy of the segmentation (A, B) as,

$$\text{Jaccard} = \max\left(\frac{J(S, A) + J(Q, B)}{2}, \frac{J(S, B) + J(Q, A)}{2}\right).$$

6.2 Sparsification algorithm for NCut

We compare our new segmentation method with the original normalized cut formulations, NCut, using the graph H described in section 3.2. In practice we sparsify the graph H to scale the eigenvector computation to large images. Again, we accomplish this using importance sampling. Let H be the graph with weights defined by Equation (3). To sample one edge from H we first select a pixel i uniformly at random. We then draw a location x from a Normal distribution centered at $X(i)$ with variance σ_x^2 and select the pixel j closest to that location. We add the edge $\{i, j\}$ to G' with weight $w'(i, j) = \exp(\|I(i) - I(j)\|^2/2\sigma_f^2)$. We repeat this process m times. In the following experiments we used $m = 100|V|$ to sparsify H .

6.3 Evaluation of NCut and MixNCut without sparsification

In order to demonstrate the efficacy of our approach compared to the traditional normalized cuts method using the graph construction given in Equation (3), we ran both NCut and MixNCut without the graph sparsification step of each algorithm. Figure 3 shows the segmentation accuracy of each method on small

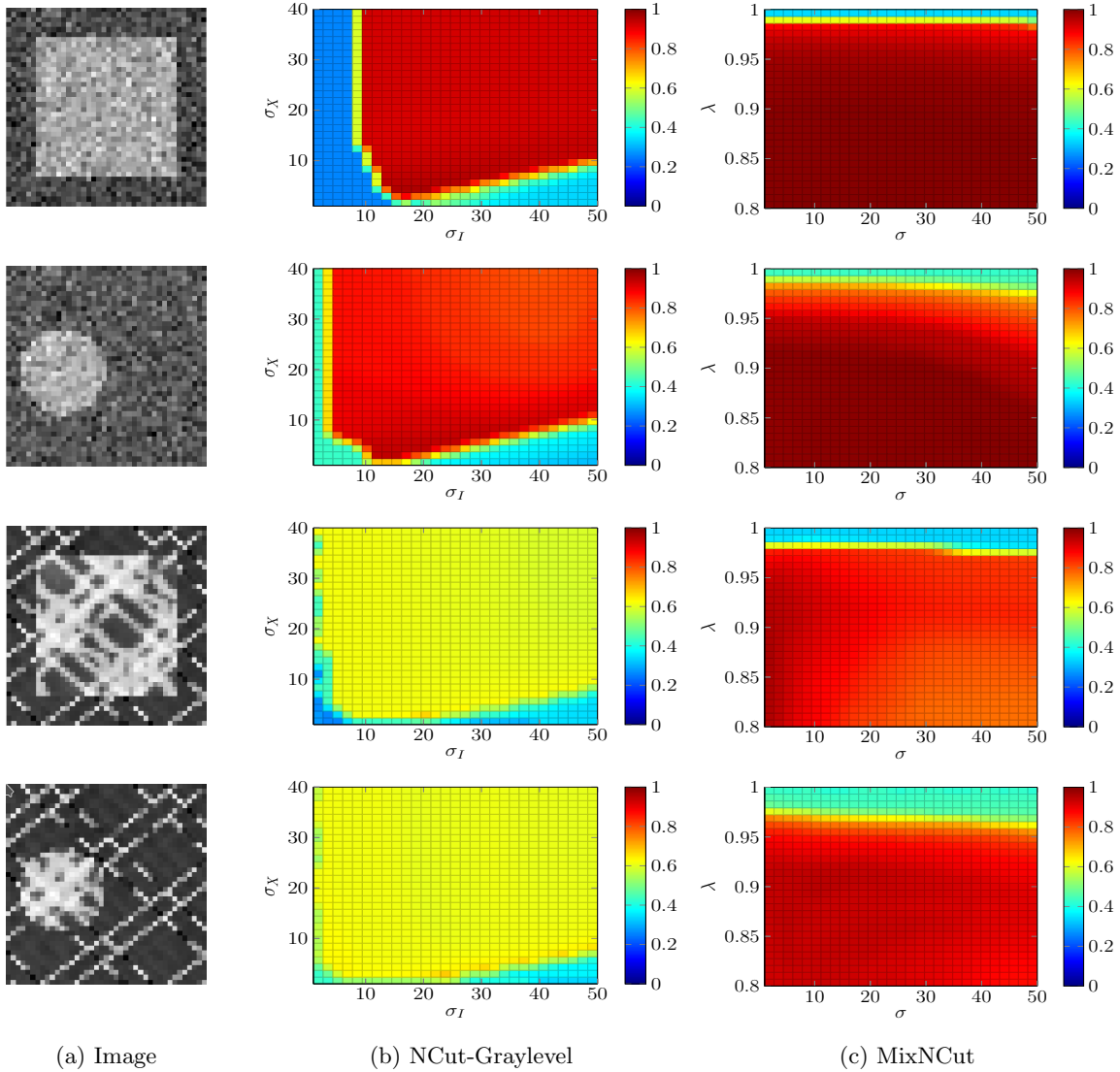


Figure 3: Evaluation of Ncut and MixNcut without graph sparsification on 40×40 images. Column (a) shows the input images. Column (b) shows the Jaccard value obtained with Ncut for a range of values for σ_I and σ_X . Column (c) shows the Jaccard value obtained with MixNcut over various combinations of σ and λ . In these experiments Ncut averaged 6.31 ± 5.28 s of processing time and MixCut, 1.35 ± 0.25 s.

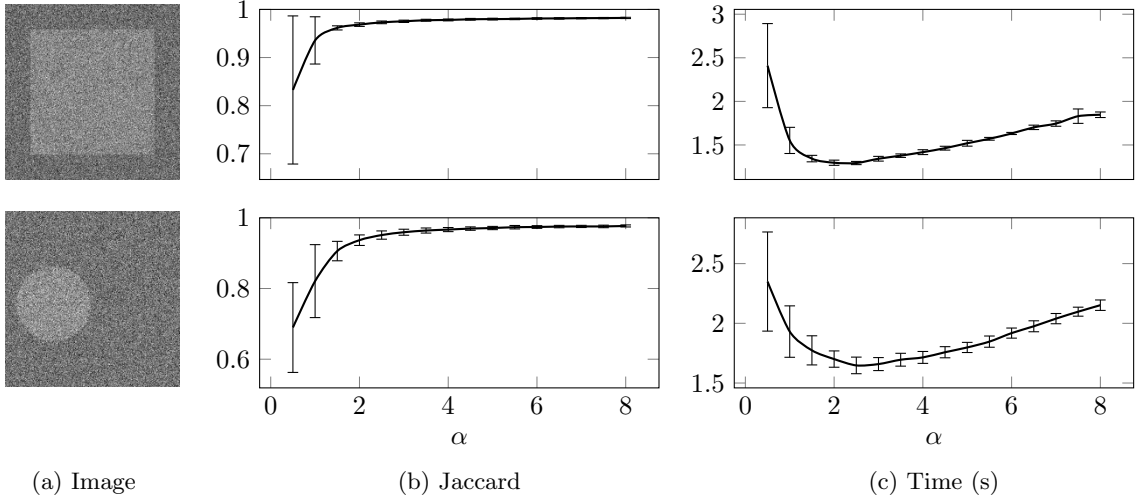


Figure 4: Impact of graph sparsification using various values for α . Column (a) shows two different test images of size 200×200 . Column (b) shows the average and standard deviation of the Jaccard accuracy measure over 100 runs of MixCut with randomized graph sparsification. Column (c) shows the average and standard deviation of the running time of the algorithm. Here, we set $\lambda = 0.95$ and $\sigma = 1$.

synthetic images under various combinations of parameters. The input images were selected to have regions of different sizes and textures.

These results demonstrate the effect of low values for either σ_I or σ_X for the normalized cut method. In this setting, the graph H is very sparse and potentially disconnected, which hinders the algorithm’s segmentation performance. We also see the performance of NCut decays in textured images.

The MixNCut method achieves good results in all scenarios, including in images with textured and/or unbalanced images.

6.4 Impact of edge sampling on MixNCut

Figure 4 shows how our algorithm performs under different values of α , where $m = \alpha|V|$ is the number of sampled edges of G_{data} when sparsified according to the algorithm in Section 5.1. Our method obtains satisfactory results in terms of Jaccard index for $\alpha \geq 2$ and achieves almost perfect segmentations for when $\alpha \geq 4$. Furthermore, our method has the lowest processing time when α is close to 2. The increase in computation time for values of α smaller than 2 is due to the slow convergence of the Lanczos algorithm in that regime. Having that in mind, the number of sampled edges used to sparsify G_{data} was set to $m = 2|V|$ for the experiments with MixNCut in the following sections.

6.5 The role of λ on MixNCut

Figure 5 depicts the impact of varying λ in our method. When λ is smaller MixCut outputs a segmentation where fine image structures are preserved. As λ increases, the resulting eigenvector is blurred, leading to segmentations with without holes or small/thin structures. These results demonstrate the role that λ , and therefore G_{grid} , plays to enforce spatial coherence within each region, whereas G_{data} promotes the fitting of the color data in each region.

6.6 Experiments in Real Images

We tested our method on real images from a variety of datasets including the Berkeley Segmentation Dataset [29], the Plant Seedlings Dataset [30], the Grabcut dataset [31], the PASCAL VOC dataset [32] and a

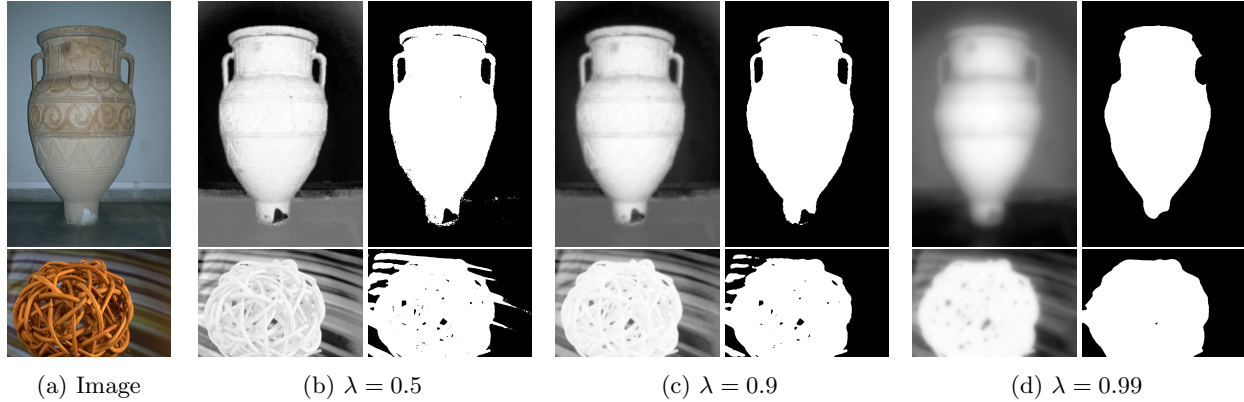


Figure 5: MixNCut results. Column (a) shows the input image. Columns (b)-(d) show the eigenvectors (on the left) and segmentations (on the right) given by MixNCut for various values of λ when $\sigma = 1$.

Scanning Electron Microscope (SEM) dataset [33]. Figure 6 shows some of the results we obtained, comparing the original normalized cuts formulation with our new approach. We can see in these examples how the new approach can segment challenging images in a variety of settings, often outperforming the original normalized cuts formulation.

Figure 7 illustrates segmentation results using MixNCut to partition an image into 3 regions. In this case we follow the approach suggested in [34] and [5], using K -means with $K = 3$ to cluster the pixels using the second and third largest eigenvector of the transition matrix P in Equation (4).

For each example in these figures, we ran the algorithms using different parameter values (specified in the next section), and show the best result among the different runs.

6.7 Experiments in Synthetic Images

For a quantitative evaluation we used images with Brodatz textures [35]. To generate input images, we mixed pairs of textures using different ground-truth segmentation patterns.

We compare our method to NCut, using either graylevel intensities or “texture features”, where we use the magnitudes of the response of 12 Gabor filters (3 wavelengths and 4 orientations) to define appearance vectors for each pixel.

Figure 8 shows some of the test images in our dataset along with the computed eigenvectors and segmentations arising with the proposed MixNCut method and the different NCut formulations. The poor segmentation performance of NCut defined over the raw pixel values (Figure 8b) can be attributed to its inability to handle the complex appearance of textured regions. This issue is partially solved when texture features defined by Gabor filters are considered, but it has the drawback of over-smoothing region boundaries (first and third examples in Figure 8c). In fact, in some extreme cases, it misses an entire small region. On the other hand, the new MixNCut method defined directly in terms of raw pixel values finds near optimal segmentations in all of these examples, preserving well the region boundaries and outperforming both baselines. This is due to its capacity to model long range relationships without relying on filtering methods.

We also compare our proposed algorithm to several state of art texture segmentation methods. They include Level Set segmentation using Wasserstein Distances (LSWD) [36], Factorization Based Segmentation (FBS) [37], Projective Non-Negative Matrix Factorization on a Graph (PNMF) [38] and ORTSEG [39]. In all these methods we try different combinations of parameters and select the ones that performed the best in our data based on their Jaccard value. We run our comparison using the implementations provided by their authors. FBS and PNMf use a bank of Gabor filters to define texture features, whereas LSWD and ORTSEG use of local image histograms.

Furthermore, we compare our method with the Multi-view Spectral Clustering (MVSC) algorithm [16]. MVSC uses a convex combination of two Laplacian matrices representing different graphs/views. For MVSC,

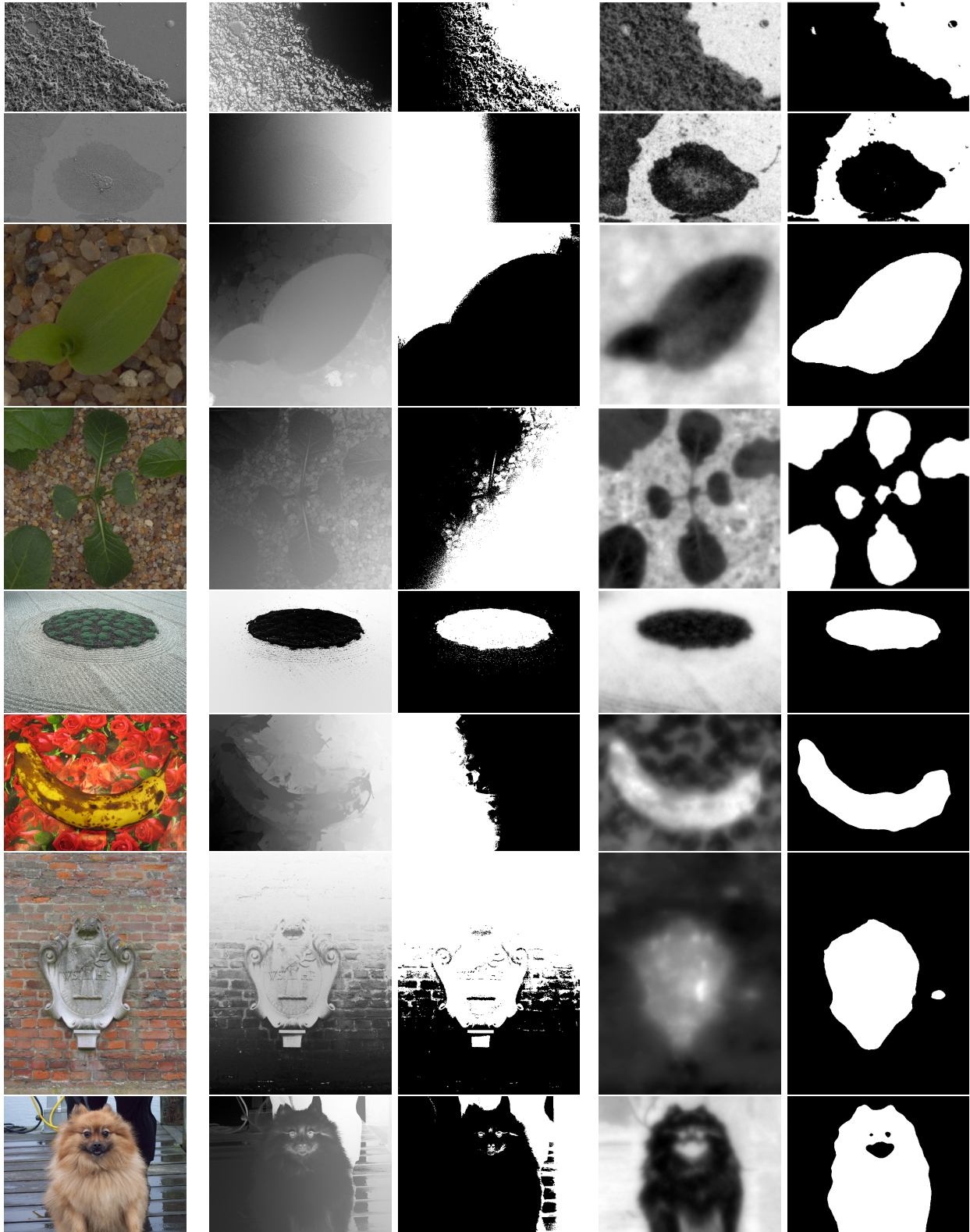


Figure 6: Segmentation results. Column (a) shows the input images. Column (b) shows the eigenvector found by the original NCut formulation on the left and the segmentation result on the right. Column (c) shows the eigenvector and segmentation found by the MixNCut method.

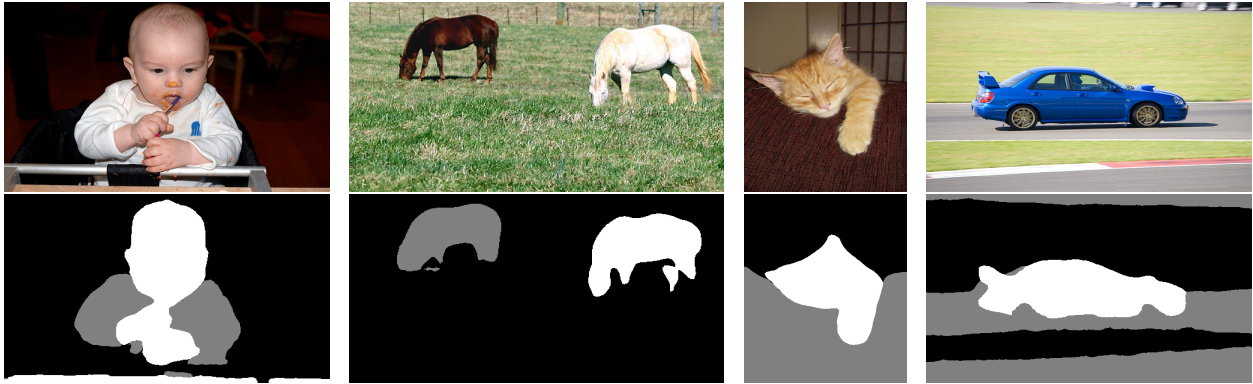


Figure 7: Segmentation results using the proposed method for images with more than 2 regions.

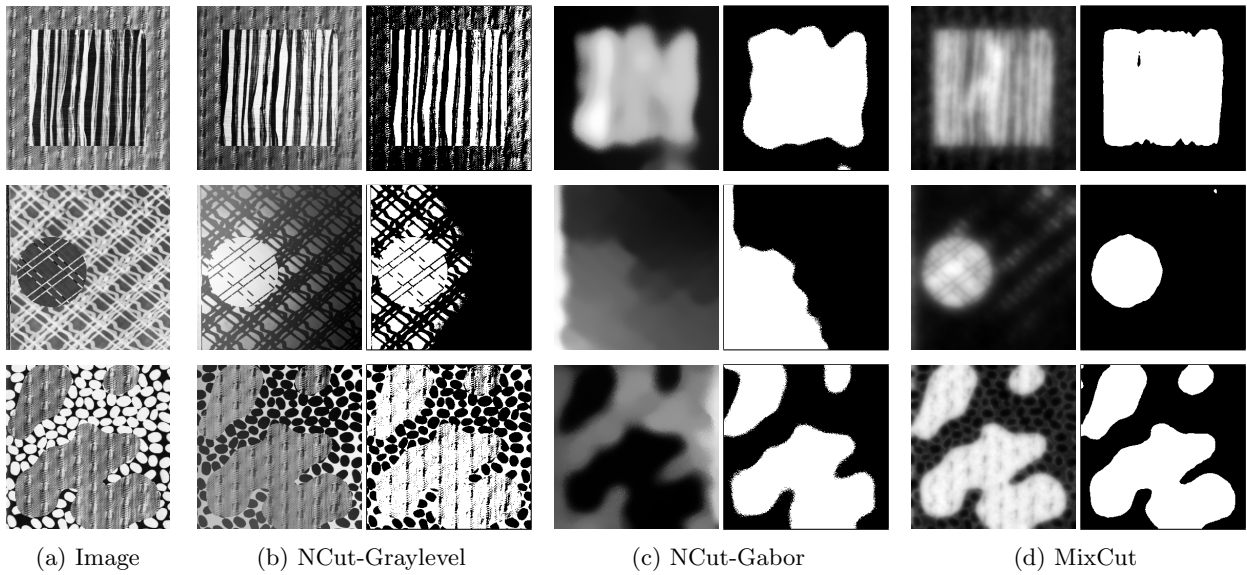


Figure 8: Comparing Ncut-Graylevel, Ncut-Gabor, and MixNCut on textured images. Column (a) shows the input images. Column (b) shows the eigenvector found by the original Ncut formulation on the left and the segmentation result on the right. Column (c) shows the eigenvector found by Ncut with Gabor features on the left and the segmentation result on the right. Column (d) shows the eigenvector found by the new MixNCut formulation on the left and the segmentation result on the right.

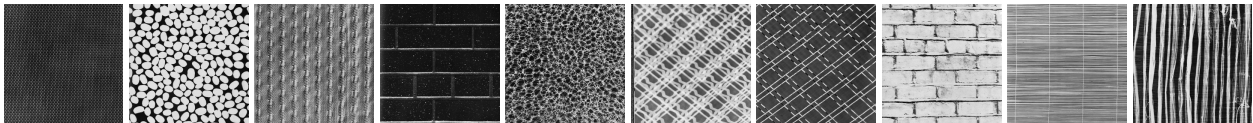

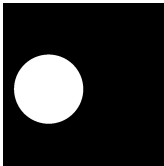
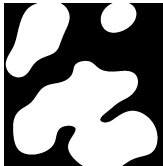


Figure 9: Brodatz Patterns used in the synthetic experiments

Table 1: Evaluation of different segmentation methods on textured images. The table summarizes accuracy and running time of each method on images with different ground-truth segmentations.

Method	Jaccard value			Time (s)
				
MixNCut	0.906 ± 0.080	0.876 ± 0.107	0.842 ± 0.133	11.27
NCut-Graylevel	0.541 ± 0.128	0.470 ± 0.126	0.538 ± 0.130	9.53
NCut-Gabor	0.800 ± 0.173	0.779 ± 0.197	0.661 ± 0.193	13.14
MVSC	0.731 ± 0.273	0.766 ± 0.214	0.656 ± 0.284	11.19
LSWD	0.852 ± 0.129	0.794 ± 0.198	0.828 ± 0.119	67.08
ORTSEG	0.853 ± 0.149	0.668 ± 0.274	0.826 ± 0.135	1.56
FBS	0.850 ± 0.151	0.778 ± 0.214	0.804 ± 0.146	0.10
PNMF	0.907 ± 0.092	0.733 ± 0.142	0.857 ± 0.085	14.59

we input the weight matrices of G_{grid} and G_{data} to their algorithm and tested the results using the same choices for λ and σ as in MixCut, picking the best segmentation according to the Jaccard value.

For the quantitative experiments, we use all pairings of the 10 textures in Figure 9 with three different ground-truth segmentations shown in Table 1 to generate three sets of images. We compute the mean accuracy of each method on each set of images using several parameter combinations. Table 1 summarizes the best mean accuracy obtained with each method on each set of inputs. The table also shows the average running time of each method. We see the new MixNCut approach obtains near perfect accuracy (Jaccard ≈ 1) on all ground-truth patterns, outperforming the other methods, specially the spectral ones.

7 Conclusion

We introduced a new spectral method for image segmentation that can segment challenging images while working directly with raw pixel values, without any pre-processing or filtering. The approach is based on a novel combination of appearance and spatial grouping cues using two different graphs. We use a dense graph to capture the appearance of regions. This leads to non-parametric models of region appearance. We also introduced a technique that can be used to sparsify the resulting graph to ease the computational burden of spectral segmentation. Our results show that long range pairwise interactions can capture the appearance of textured regions and significantly improve the performance of graph-based segmentation methods. The proposed method is practical and it can be applied to different types of images (natural scenes, biomedical, textures, etc.) that arise in a variety of application.

Broader Impact

Image segmentation has a variety of applications that can benefit all of society. For example, segmentation methods may enable advances in biomedical image analysis (including for medical diagnosis and treatment), tele-conferencing technology, human-computer-interaction, remote sensing, and robotics. However, there are also potential uses with questionable ethics, including mass surveillance and military applications.

References

- [1] J. Shi and J. Malik, “Normalized cuts and image segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888–905, 2000.
- [2] W. Donath and A. Hoffman, “Lower bounds for the partitioning of graphs,” *IBM Journal of Research and Development*, vol. 17, no. 5, pp. 420–425, 1973.
- [3] M. Fiedler, “Algebraic connectivity of graphs,” *Czechoslovak mathematical journal*, vol. 23, no. 2, pp. 298–305, 1973.
- [4] G. Cheung, E. Magli, Y. Tanaka, and M. K. Ng, “Graph spectral image processing,” *Proceedings of the IEEE*, vol. 106, no. 5, pp. 907–930, 2018.
- [5] M. Meila and J. Shi, “Learning segmentation by random walks,” in *Advances in Neural Information Processing Systems*, pp. 873–879, 2001.
- [6] S. Maji, N. K. Vishnoi, and J. Malik, “Biased normalized cuts,” in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2057–2064, IEEE, 2011.
- [7] S. E. Chew and N. D. Cahill, “Semi-supervised normalized cuts for image segmentation,” in *IEEE International Conference on Computer Vision*, pp. 1716–1723, 2015.
- [8] P. Perona and W. Freeman, “A factorization approach to grouping,” in *European Conference on Computer Vision*, pp. 655–670, Springer, 1998.
- [9] C. Fowlkes, S. Belongie, F. Chung, and J. Malik, “Spectral grouping using the nystrom method,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 26, no. 2, pp. 214–225, 2004.
- [10] F. Lin and W. W. Cohen, “Power iteration clustering,” in *International Conference on Machine Learning*, 2010.
- [11] J. Malik, S. Belongie, T. Leung, and J. Shi, “Contour and texture analysis for image segmentation,” *International Journal of Computer Vision*, vol. 43, no. 1, pp. 7–27, 2001.
- [12] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, “Contour detection and hierarchical image segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 5, pp. 898–916, 2010.
- [13] S. Bickel and T. Scheffer, “Multi-view clustering,” in *ICDM*, pp. 19–26, 2004.
- [14] V. R. De Sa, “Spectral clustering with two views,” in *ICML Eorkshop on Learning with Multiple Views*, pp. 20–27, 2005.
- [15] A. Kumar, P. Rai, and H. Daume, “Co-regularized multi-view spectral clustering,” in *Advances in Neural Information Processing Systems*, pp. 1413–1421, 2011.
- [16] D. Zhou and C. J. Burges, “Spectral clustering and transductive learning with multiple views,” in *International Conference on Machine learning*, pp. 1159–1166, 2007.
- [17] L. Zong, X. Zhang, X. Liu, and H. Yu, “Weighted multi-view spectral clustering based on spectral perturbation,” in *AAAI Conference on Artificial Intelligence*, 2018.
- [18] P. Mercado, A. Gautier, F. Tudisco, and M. Hein, “The power mean laplacian for multilayer graph clustering,” in *International Conference on Artificial Intelligence and Statistics*, vol. 84, pp. 1828–1838, 2018.
- [19] X. Xia and B. Kulis, “W-net: A deep model for fully unsupervised image segmentation,” *arXiv:1711.08506*, 2017.

- [20] M. Tang, F. Perazzi, A. Djelouah, I. Ben Ayed, C. Schroers, and Y. Boykov, “On regularized losses for weakly-supervised cnn segmentation,” in *European Conference on Computer Vision*, pp. 507–522, 2018.
- [21] S. Ghosh, N. Das, I. Das, and U. Maulik, “Understanding deep learning techniques for image segmentation,” *ACM Computing Surveys*, vol. 52, no. 4, pp. 1–35, 2019.
- [22] C. Tomasi and R. Manduchi, “Bilateral filtering for gray and color images,” in *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*, pp. 839–846, IEEE, 1998.
- [23] Y. Boykov and V. Kolmogorov, “Computing geodesics and minimal surfaces via graph cuts,” in *IEEE International Conference on Computer Vision*, 2003.
- [24] D. Mumford and J. Shah, “Optimal approximations by piecewise smooth functions and associated variational problems,” *Communications on pure and applied mathematics*, vol. 42, no. 5, pp. 577–685, 1989.
- [25] R. Jenssen, D. Erdogmus, J. Principe, and T. Eltoft, “The laplacian PDF distance: A cost function for clustering in a kernel feature space,” in *Advances in Neural Information Processing Systems*, pp. 625–632, 2005.
- [26] M. Tang, L. Gorelick, O. Veksler, and Y. Boykov, “Grabcut in one cut,” in *IEEE International Conference on Computer Vision*, pp. 1769–1776, 2013.
- [27] D. P. Williamson and D. B. Shmoys, *The design of approximation algorithms*. Cambridge University press, 2011.
- [28] P. Jaccard, “Étude comparative de la distribution florale dans une portion des alpes et des jura,” *Bull Soc Vaudoise Sci Nat*, vol. 37, pp. 547–579, 1901.
- [29] D. Martin, C. Fowlkes, D. Tal, and J. Malik, “A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics,” in *IEEE International Conference on Computer Vision*, vol. 2, pp. 416–423, July 2001.
- [30] T. M. Giselsson, R. N. Jørgensen, P. K. Jensen, M. Dyrmann, and H. S. Midtiby, “A public image database for benchmark of plant seedling classification algorithms,” *arXiv preprint arXiv:1711.05458*, 2017.
- [31] C. Rother, V. Kolmogorov, and A. Blake, “Grabcut” interactive foreground extraction using iterated graph cuts,” *ACM transactions on graphics (TOG)*, vol. 23, no. 3, pp. 309–314, 2004.
- [32] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes challenge: A retrospective,” *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.
- [33] R. Aversa, M. H. Modarres, S. Cozzini, and R. Ciancio, “NFFA-EUROPE - SEM dataset,” 2018.
- [34] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *Advances in Neural Information Processing Systems*, pp. 849–856, 2002.
- [35] P. Brodatz, *Textures: a photographic album for artists and designers*. Dover Pubns, 1966.
- [36] K. Ni, X. Bresson, T. Chan, and S. Esedoglu, “Local histogram based segmentation using the wasserstein distance,” *International Journal of Computer Vision*, vol. 84, no. 1, pp. 97–111, 2009.
- [37] J. Yuan, D. Wang, and A. M. Cheriyyadat, “Factorization-based texture segmentation,” *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 3488–3497, 2015.

- [38] C. G. Bampis, P. Maragos, and A. C. Bovik, “Projective non-negative matrix factorization for unsupervised graph clustering,” in *IEEE International Conference on Image Processing*, pp. 1255–1258, IEEE, 2016.
- [39] M. T. McCann, D. G. Mixon, M. C. Fickus, C. A. Castro, J. A. Ozolek, and J. Kovacević, “Images as occlusions of textures: A framework for segmentation,” *IEEE transactions on image processing*, vol. 23, no. 5, pp. 2033–2046, 2014.