

**UCLA**

**UCLA Previously Published Works**

**Title**

Diffuse Interface Models on Graphs for Classification of High Dimensional Data

**Permalink**

<https://escholarship.org/uc/item/4hs5472g>

**Journal**

SIAM REVIEW, 58(2)

**ISSN**

0036-1445

**Authors**

Bertozzi, Andrea L  
Flenner, Arjuna

**Publication Date**

2016

**DOI**

10.1137/16M1070426

Peer reviewed

# Diffuse Interface Models on Graphs for Classification of High Dimensional Data\*

Andrea L. Bertozzi<sup>†</sup>

Arjuna Flenner<sup>‡</sup>

**Abstract.** This paper is a republication of an MMS paper [A. L. Bertozzi and A. Flenner, *Multiscale Model. Simul.*, 10 (2012), pp. 1090–1118] describing a new class of algorithms for classification of high dimensional data. These methods combine ideas from spectral methods on graphs with nonlinear edge/region detection methods traditionally used in the PDE-based imaging community. The algorithms use the Ginzburg–Landau functional, extended to a graphical framework, which has classical PDE connections to total variation minimization. Convex splitting algorithms allow us to quickly find minimizers of the proposed model and take advantage of fast spectral solvers of linear graph-theoretic problems. We review the diverse computational examples presented in the original paper, involving both basic clustering and semisupervised learning for different applications. Case studies include feature identification in images, segmentation in social networks, and segmentation of shapes in high dimensional datasets. Since the paper’s publication four years ago a body of research has developed centering on this class of methods from a variety of perspectives, including analytical results involving convergence and connection to graph cuts, extensions of the method using the MBO scheme and to problems involving social networks, and comparison to other approaches. We review these more recent results in an extended conclusion section of the paper.

**Key words.** Nyström extension, diffuse interfaces, image processing, high dimensional data

**AMS subject classifications.** 68U10, 49-04, 49M25, 62H30, 68T10, 91C20

**DOI.** 10.1137/16M1070426

**Introduction.** The work develops a new algorithm for classification of high dimensional data on graphs. The method bridges gaps between the PDE-image processing community, methods related to  $L^1$  compressive sensing, and the graph-theoretic community. The algorithm is inspired by diffuse interface methods that have long been used in physical science modeling of problems with free boundaries and inter-

---

\*Published electronically May 5, 2016. This paper originally appeared in *Multiscale Modeling & Simulation*, Volume 10, Number 3, 2012, pages 1090–1118. The original MMS paper was supported by ONR grants N000140810363, N000141010221, N000141210040, and N0001411AF00002; NSF grants DMS-0914856 and DMS-1118971; and AFOSR MURI grant FA9550-10-1-0569. The SIGEST version of the paper was additionally supported by NSF grant DMS-1417674 and ONR grant N00014-16-1-2119. This work was performed by an employee of the U.S. Government or under U.S. Government contract. The U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or allow others to do so, for U.S. Government purposes. Copyright is owned by SIAM to the extent not limited by these rights.

<http://www.siam.org/journals/sirev/58-2/M107042.html>

<sup>†</sup>Department of Mathematics, University of California, Los Angeles, CA 90095 (bertozzi@math.ucla.edu).

<sup>‡</sup>Physics and Computational Sciences, Naval Air Weapons Center, China Lake, CA (arjuna.flenner@navy.mil).

faces. These same ideas also arise in compressive sensing because of the connection between the diffuse interface energies and total variation (TV) minimization. Thus it is natural to extend these concepts beyond classical continuum modeling to the discrete graph framework. By implementing various techniques for fast linear algebra and eigenfunction decomposition, we can design methods that are quite good in terms of performance and speed. They are also versatile, and we illustrate examples on a variety of datasets, including congressional voting records, high dimensional test data, and machine learning in image processing. We note that there is a large literature (discussed later in the paper) connecting graph cuts to spectral properties of the graph Laplacian—most notably work on spectral clustering and also on Cheeger cuts and the Cheeger ratio. This paper introduces a completely different approach to approximating graph cuts by considering a nonlinear approximation built on curvature driven models from PDE. This approximation, involving the Ginzburg–Landau functional, can be efficiently minimized using spectral solvers, thus taking advantage of modern machine learning algorithms for graph Laplacians.

This paper is structured as follows: In section 1 we review diffuse interface methods in Euclidean space and convex splitting methods for minimization. These well-known constructions make heavy use of the classical Laplace operator, and our new algorithms involve extensions of this idea to a more general graph Laplacian. Section 2 reviews some of the notation and definitions of the graph Laplacian, and this discussion contains a level of detail appropriate for readers less familiar with this machinery. Included in this section are a review of segmentation using spatial clustering and a discussion of various normalization conventions for these linear operators on graphs, in connection to real world problems such as machine learning in image analysis. Section 3 presents the computational algorithm as motivated by classical pseudospectral methods for nonlinear PDE. Section 4 discusses computational examples, while section 5 connects our work to other work on graph cuts. Section 6 reviews the relevant literature since the publication of the original MMS version of this paper.

**1. Background on Diffuse Interfaces, Image Processing, and Convex Splitting Methods.** Diffuse interface models in Euclidean space are often built around the Ginzburg–Landau (GL) functional

$$(1.1) \quad GL(u) = \frac{\epsilon}{2} \int |\nabla u|^2 dx + \frac{1}{\epsilon} \int W(u) dx,$$

where  $W$  is a double well potential. For example,  $W(u) = \frac{1}{4}(u^2 - 1)^2$  has minimizers at plus and minus one. The operator  $\nabla$  denotes the spatial gradient operator, and the first term in  $GL$  is  $\epsilon/2$  times the  $H^1$  seminorm of  $u$ . The small parameter  $\epsilon$  represents a spatial scale, the diffuse interface scale.

The model is called “diffuse interface” because there is a competition between the two terms in the energy functional. Upon minimization of this functional, the double well will force  $u$  to go to either one or minus one; however, the  $H^1$  term forces  $u$  to have some smoothness, thereby removing sharp jumps between the two minima of  $W$ . The resulting minimization leads to regions where  $u$  is approximately one, regions where it is approximately minus one, and a very thin  $O(\epsilon)$  scale transition region between the two. Thus the minimizer appears to have two phases with an interface between them. For this reason such models are often referred to as “phase field” models when one considers dynamic evolution equations built around this energy functional. There are several interesting features of GL minimizers. For example, the transition region between the two phases typically has some length associated with

it, and the GL functional is roughly proportional to this length. This can be made rigorous by considering the notion of Gamma convergence of the GL functional. It is known to converge [50] to the TV seminorm,

$$GL(u) \rightarrow_{\Gamma} C|u|_{TV}.$$

Diffuse interface models with spatial gradient operators are multiscale because several length scales exist in the model, the smallest being the diffuse interface scale  $\epsilon$ . The construction of the GL functional requires  $\epsilon$  to have units of length so that the two terms have balanced units. Also, note that the double well typically restricts  $u$  to take on integral order values (between zero and one).

Multiple time scales exist in evolution equations arising from the functional (1.1). The most common examples are the Allen–Cahn equation, which is the  $L^2$  gradient descent of this functional, and the Cahn–Hilliard equation, which is a gradient descent in the  $H^{-1}$  inner product. Recent work has gone into developing efficient numerical schemes to track these dynamics [49, 82, 48, 9]; the same ideas can be used to speed up variational algorithms based on these functionals.

**1.1. The Connection between GL and Image Processing.** Perhaps the most important motivation for the method proposed in this paper is the existing literature on the use of the GL functional in image processing. We review some of that literature for the reader not directly familiar with it to better motivate the choice of methods proposed here. The GL functional is sometimes used in image processing as an alternative or a relative to the TV seminorm. Because of the Gamma convergence, these two functionals can sometimes be interchanged. Moreover, the highest order term in the GL functional is purely quadratic, allowing for fast minimization schemes in some problems. Recent advances in TV minimization procedures, e.g., split Bregman and graph cut methods [40, 19], have made this less necessary; nevertheless there are cases where the pure TV case is not enough and the diffuse interface version may be a simpler method.

One example of the GL functional in image processing is the motivation for the Esedoğlu–Tsai threshold method [28] for Chan–Vese segmentation [16]. The construction of their method is directly built on the GL functional, rather than the TV method of the original Chan and Vese paper [16]. We also note that this method was inspired by the original Merriman, Bence, and Osher (MBO) paper [57], and section 6 discusses the connection between the original MBO paper and our graph-based method in more detail. Another example is work by March (see [26, 14]) and Esedoglu (see [25, 27]). The work by Ambrosio and Tortorelli [3] is well known in image processing for diffuse interface approximations.

In a typical application we want to minimize an energy functional of the form

$$E(u) = GL(u) + \lambda F(u, u_0),$$

where  $F(u, u_0)$  is a fitting term to known data. In the case of denoising,  $F(u, u_0)$  is often just an  $L^2$  fit,  $\int (u - u_0)^2$ . In the case of deblurring it is  $\int (K * u - u_0)^2$ , or the  $L^2$  of the blurred solution with the data. For inpainting we often have an  $L^2$  fit to known data in the region where the data is known, i.e.,  $\int_{\Omega} (u - u_0)^2$ . In some instances in the above a different norm is used, e.g.,  $L^1$  or other norms. In the case of Cahn–Hilliard-based inpainting, the method is not strictly a gradient flow [8, 7], but rather is based on gradient flows. In fact, the method is a sort of hybrid in which the  $L^2$  least-squares fitting term is paired with the Cahn–Hilliard  $H^{-1}$  dynamics. The

result is a method that achieves both required boundary conditions for inpainting, namely, continuation of both grayscale information and direction of edges across the inpainting domain. The higher order evolution is important for that application and is related to the geometry of the problem. The Cahn–Hilliard inpainting problem is a good example of a method where new fast algorithms for TV minimization have yet to be developed due, in part, to the unusual boundary conditions. The convex splitting method (see below) allows for very fast inpainting using a Cahn–Hilliard method and the GL functional.

The energy  $E(u)$  can be minimized in the  $L^2$  sense using a gradient descent, which gives us a modified Allen–Cahn equation

$$u_t = -\frac{\delta GL}{\delta u} - \lambda \frac{\delta F}{\delta u} = \epsilon \Delta u - \frac{1}{\epsilon} W'(u) - \lambda \frac{\delta F}{\delta u}.$$

This can be evolved to steady state to obtain a local minimizer of the energy  $E$ . We note that, in general, especially for the GL functional,  $E$  is not convex and thus may have multiple local energy minima. The result is that the long time behavior of the solution of the modified Allen–Cahn equation may depend on the initial condition.

**1.2. Convex Splitting and Time Stepping of the GL Functional.** One of the reasons to choose the GL functional instead of TV is that the minimization procedure for GL often involves the first variation of GL for which the highest order term, involving the Laplace operator, is linear. Thus if one has fast solvers for the Laplace operator or relatives of it, one can take advantage of this in designing convex splitting schemes discussed below. This class of schemes will be developed in the context of graph-based methods later in this paper. To motivate the ideas for the reader unfamiliar with the technique, we provide some background on this class of methods in relation to work involving the GL functional.

A particular class of fast solvers is one in which the Laplacian can be transformed so that the operator diagonalizes. A classical example would be the fast Fourier transform (FFT), which transforms the Laplace operator to multiplication by  $-|k|^2$ , where  $k$  is the wave number of the Fourier mode. The FFT works because the Fourier modes are also eigenfunctions of the Laplace operator. An example of this use in long time solutions of the Cahn–Hilliard equation is discussed in detail in [82]. Other recent advances for fast Poisson solvers could be used as well (see, e.g., [56]). In our graph-based examples we use fast methods for directly diagonalizing the graph Laplacian, either through standard sparse linear algebra routines or, in the case of fully connected weighted graphs, Nyström extension methods.

Convex splitting schemes are based on the idea that an energy functional can be written as the sum of convex and concave parts,

$$E(u) = E_{\text{vex}}(u) - E_{\text{cave}}(u),$$

where this decomposition is certainly not unique because we can add and subtract any convex function and not change  $E$  but certainly change the convex/concave splitting. The idea behind convex splitting for the gradient descent problem is to perform a time-stepping scheme in which the convex part is done implicitly and the concave part explicitly. More precisely, the convex splitting scheme is

$$(1.2) \quad \frac{u^{n+1} - u^n}{dt} = -\frac{\delta E_{\text{vex}}}{\delta u}(u^{n+1}) + \frac{\delta E_{\text{cave}}}{\delta u}(u^n).$$

The art then lies in choosing the splitting so that the resulting scheme is stable and also computationally efficient to solve. This method was popularized by a well-known but unpublished manuscript by Eyre [29]. It has been used to solve the Cahn–Hilliard equation on large domains and on long time intervals [82] and also in imaging applications involving Cahn–Hilliard [8] and a wavelet version of the method proposed here for graphs [21]. This same idea has also been directly discussed in the context of general minimization procedures for nonconvex functionals [89].

**2. Generalizations of the GL Functional to Graphs.** One can consider a generalization of the GL functional to graphs. This is in the same spirit as the work of Dobrosotskaya and the first author [21] generalizing the GL functional to wavelets. In their work they construct a linear operator with features similar to those of the Laplace operator; however, the eigenfunctions are the wavelet basis for some choice of wavelets. The natural choice of eigenvalues are ones that scale like the inverse square of the length scale of the wavelet basis functions, much in the same way that the eigenvalues of the Laplace operator are the inverse square of the period of the corresponding eigenfunction.

In this section we describe how to generalize the GL functional, or, more precisely, its  $L^2$  gradient flow, to the case of functions defined on graphs [17]. One challenge is the normalization of the Laplacian due to the fact that we are working with purely discrete functionals that may not have a direct spatial embedding. We include some details for the reader not familiar with the graph-based literature. This section is organized into five subsections discussing how the graph Laplacian can be incorporated into a GL-type functional. First, we introduce the graph Laplacian using weight functions. We specifically introduce the notation used in this paper. In order to provide insight into our algorithm, we continue in the next subsection with a short discussion on previous segmentation algorithms using the graph Laplacian. Next, the benefits of the normalized graph Laplacian are discussed. The graph-based GL functional for segmentation is then introduced, and the last subsection mentions techniques for creating weight functions and thus building interesting graph Laplacians.

**2.1. Graph Definitions and Notation.** Consider an undirected graph  $G = (V, E)$  with vertex set  $V = \{\nu_n\}_{n=1}^N$  and edge set  $E$ . A weighted undirected graph [17] has an associated weight function  $w : V \times V \rightarrow R$  satisfying  $w(\nu, \mu) = w(\mu, \nu)$  and  $w(\nu, \mu) \geq 0$ . The degree of a vertex  $\nu \in V$  is defined as

$$(2.1) \quad d(\nu) = \sum_{\mu \in V} w(\nu, \mu).$$

The degree matrix  $D$  can then be defined as the  $N \times N$  diagonal matrix with diagonal elements  $d(\nu)$ .

The size of a subset  $A \subset V$  will be important for segmentation using graph theory, and there are two important size measurements. For  $A \subset V$  define

$$(2.2) \quad |A| := \text{the number of vertices in } A,$$

$$(2.3) \quad \text{vol}(A) := \sum_{\nu \in A} d(\nu).$$

The topology of the graph also plays a role. A subset  $A \subset V$  of a graph is connected if any two vertices in  $A$  can be joined by a path such that all the points also lie in  $A$ . A subset of  $A$  is called a connected component if it is connected and if  $A$  and  $\bar{A}$  are

not connected. The sets  $A_1, A_2, \dots, A_k$  form a partition of the graph if  $A_i \cap A_j = \emptyset$  and  $\cup_k A_k = V$ .

The graph Laplacian is the main tool for graph theory-based segmentation. Define the graph Laplacian  $L(\nu, \mu)$  as

$$(2.4) \quad L(\nu, \mu) = \begin{cases} d(\nu) & \text{if } \nu = \mu, \\ -w(\nu, \mu) & \text{otherwise.} \end{cases}$$

The graph Laplacian can be written in matrix form as  $L = D - W$ , where  $W$  is the matrix  $w(\nu, \mu)$ . The following definition and property of  $L$  are important:

1. (quadratic form) For every vector  $u \in R^N$

$$(2.5) \quad \langle u, Lu \rangle = \frac{1}{2} \sum_{\mu, \nu \in V} w(\nu, \mu) (u(\nu) - u(\mu))^2.$$

2. (eigenvalue)  $L$  has  $N$  nonnegative, real valued eigenvalues with  $0 = \tilde{\lambda}_1 \leq \tilde{\lambda}_2 \leq \dots \leq \tilde{\lambda}_N$ , and the eigenvector of  $\tilde{\lambda}_1$  is the constant  $N$ -dimensional one vector  $1_N$ .

The quadratic form is exploited to define a minimization procedure as in the Allen-Cahn equation above. The eigenvalue condition gives limitations on the spectral decomposition of the matrix  $L$ . These spectral properties are essential for the spectral clustering algorithms discussed below.

There are two popular normalization procedures for the graph Laplacian, with consequences for segmentation [17, 83]. The normalization primarily used here is the symmetric Laplacian  $L_s$  defined as

$$(2.6) \quad L_s = D^{-1/2} L D^{-1/2} = I - D^{-1/2} W D^{-1/2}.$$

The symmetric Laplacian is named as such since it is a symmetric matrix; it is particularly useful for the nonlocal means graph for images (see Figure 2.1). The random walk Laplacian is another important normalization given by

$$(2.7) \quad L_w = D^{-1} L = I - D^{-1} W.$$

The random walk Laplacian is closely related to discrete Markov processes, and we discuss the use of the random walk Laplacian in section 5.2. The papers [60, 80] provide additional discussions of the different versions of  $L$  in the context of curvature-driven methods for graph cuts. We note that the unnormalized Laplacian is used to prove Gamma convergence of the graph-GL functional to the graph TV functional in [79].

The spectrum of the graph Laplacian is used for graph segmentation, and some well-known results are collected here for future reference. The spectra of  $L_s$  and  $L_w$  are the same, but the eigenvectors are different. The easily verifiable spectral relationships between  $L_w$  and  $L_s$  follow:

1.  $\tilde{\lambda}$  is an eigenvalue of  $L_w$  if and only if  $\tilde{\lambda}$  is an eigenvalue of  $L_s$ .
2.  $\psi$  is an eigenvector of  $L_w$  if and only if  $D^{1/2}\psi$  is an eigenvector of  $L_s$ .
3.  $\lambda$  is an eigenvalue of  $L_w$  with eigenvector  $\psi$  if and only if  $L\psi = \tilde{\lambda}D\psi$ .

**2.2. Choice of Similarity Function.** Here we discuss how to construct the weight functions  $w(x, y)$ , sometimes referred to as similarity functions, for specific applications involving high dimensional data. There are two factors to consider when choos-

ing  $w(x, y)$ . First, the choice of weight function must reflect the desired outcome. For segmentation, this typically involves choosing an appropriate metric on a vector space. Our examples below used the standard Euclidean norm, but other norms may be more appropriate. For example, the angle norm may work better for segmentation of hyperspectral images. A second consideration is algorithm speed. The segmentation algorithms below requires the diagonalization of  $w(x, y)$ , and this step is often the rate limiting procedure. There are two main methods for obtaining speed in the diagonalization. The first method is to use the Nyström extension described in section 3.2. This method does not require a modification of  $w(x, y)$ , and calculations on large graphs with connections between every vertex are possible.

The second method is to create a sparse graph. A sparse graph can be created by keeping only the  $N$  largest values of  $w(x, y)$  for each fixed  $x$ . Note that such a graph is not symmetric, but it can easily be made symmetric to aid in computation.

We list the two techniques to create the similarity function  $w(x, y)$  used in this paper.

1. The Gaussian function

$$(2.8) \quad w(x, y) = \exp(-\|x - y\|^2/\tau)$$

is a common similarity function. Depending on the choice of metric, this similarity function includes the Yaroslavsky filter [88] and the nonlocal means filter [12].

2. Zelnik-Manor and Perona introduced local scaling weights for sparse matrix computations [90]. They start with a metric  $d(x_i, x_j)$  between each sample point. The idea is to define a local parameter  $\sqrt{\tau(x_i)}$  for each  $x_i$ . The choice in [90] is  $\sqrt{\tau(x_i)} = d(x_i, x_M)$ , where  $x_M$  is the  $M$ th closest vector to  $x_i$ . In [90],  $M = 7$ , while in this work and [74],  $M = 10$ . The similarity matrix is then defined as

$$(2.9) \quad w(x, y) = \exp\left(-\frac{d(x, y)^2}{\sqrt{\tau(x)\tau(y)}}\right).$$

This similarity matrix is better at segmentation when there are multiple scales that need to be segmented simultaneously.

**2.3. Segmentation, Spectral Clustering, and the Graph Laplacian.** In this subsection we review some of the previous literature on spectral clustering, in which one directly uses the spectrum of the graph Laplacian to segment data. The goal of graph clustering is to partition the vertices into groups according to their similarities. Consider the weight function as a measure of the similarities; then the graph problem is equivalent to finding a partition of the vertices such that the sum of the edge weights between the groups is small compared with the sum of the edges within the groups. The weighted graph minimization algorithms in their original form are NP complete problems [83]; therefore a relaxed problem was formulated by Shi and Malik [70], where the domain of the energy functional is allowed to be real valued and such minimization problems are equivalent to spectral clustering methods.

The segmentation problem naturally generates a graph structure from a set of vertices  $v_i$ , each of which is assigned a vector  $z_i \in \mathbb{R}^K$ . For example, when considering voting records of the U.S. House of Representatives, each representative defines a vertex, and his or her voting record defines a vector. A different example arises when considering similarity between regions in image data. Each pixel defines a vertex,



and one can assign a high dimensional vector to that pixel by comparing similarities between the neighborhood around that pixel and that of any other pixel. Given such an association, a symmetric weight matrix can be created using a symmetric function  $\hat{w}(x, y) : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}_+$ . In particular, if  $\nu_i(y) = z_i$  represents the vector associated with the vertex  $\nu_i$ , then the weight matrix  $w(\nu_i, \mu_j) = \hat{w}(\nu_i(z), \mu_j(z)) = \hat{w}(z_i, z_j)$  is a positive symmetric function. We will abuse notation and not distinguish between these two functions and write  $w(\nu_i, \mu_j) = \hat{w}(z_i, z_j) = w(z_i, z_j)$ . Similar statements are true for any function  $u : V \rightarrow \mathbb{R}$ . Spectral clustering algorithms for binary segmentation consist of the following steps:

**Input:** A set of vertices  $V$  with the associated set of vectors  $Z \subset \mathbb{R}^K$ , a similarity measure  $w(x, y) : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}_+$ , and the integer  $k$  of clusters to construct.

1. Calculate the weight function  $w(x, y)$  for all  $x, y \in Z$ .
2. Compute the graph Laplacian  $L$ .
3. Compute the second eigenvector  $\psi_2$  of  $L$  or the second eigenvector  $\psi_2$  of the generalized eigenvalue problem  $L\psi = \lambda D\psi$ .
4. Segment  $\psi_2$  into two clusters using  $k$ -means (with  $k = 2$ ).

**Output:** A partition of  $V$  (or, equivalently,  $Z$ ) into two clusters  $A$  and  $\bar{A}$ .

Two characteristics of the spectral clustering algorithms should be highlighted. First, the algorithm determines clusters using a  $k$ -means algorithm. We note that the  $k$ -means algorithm is used to construct a partition of the real valued output, and any algorithm that performs this goal can be substituted for the  $k$ -means algorithm. For example, Lang [54] uses separating hyperplanes. A partitioning algorithm is needed since the relaxed problem does not force the final output function  $f$  to be binary valued. We address this problem by using the GL potential.

The second characteristic is that spectral clustering finds natural clusters through a constrained minimization problem. The constrained minimization problem exploits a finite number of eigenfunctions depending on the a priori chosen number of clusters. A significant difference in our method is that we utilize *all* the eigenfunctions in our variational problem. One can interpret this as an issue of the number of scales that need to be resolved to perform the desired classification. For spectral clustering to work, the eigenfunctions used must capture all the relevant scales in the problem. By using all the eigenfunctions, we resolve essentially all the scales in the problem, modulo the choice of  $\epsilon$ . In the classical differential equation problem,  $\epsilon$  selects a smallest length scale to be resolved for the interfacial problem. An analogous role could occur in the graph problem, and thus it would make sense to use this method on large dataset problems rather than relatively small problems, for which other methods might be simpler.

**2.4. Proper Normalization of the Graph Laplacian with Scale.** An important issue with nonlocal operators is the behavior of the operators with increased sample size. Increasing sample size for the discrete Laplace operator corresponds to decreasing the grid size, and this operator must be correctly normalized in order to guarantee convergence to the differential Laplacian. We note that in the case of the classical finite difference problem for PDEs the entire matrix is multiplied by  $N^2$ , where  $N$  is the number of vertices in one of the dimensions, which is essentially  $1/dx$ , the spatial grid size. Recall that the largest eigenvalue of the operator scales like  $N^2$  or  $1/dx^2$ , which gives a stiffness constraint for forward time stepping of the heat equation, as a function of grid size. Moreover, with this scaling, the graph Laplacian converges to the continuum differential operator in the limit of large sample size, i.e., as  $N \rightarrow \infty$ , where  $N$  is the grid resolution along one of the coordinate axes.

Proper normalization conditions for convergence also exist for the graph Laplacian. The issue of sample size also comes into play, but rather than convergence to a differential operator, we consider the density of vertices, in the case of spatial embeddings, which can be measured by the degree of each vertex. The normalized Laplace operator as defined in (2.6) is known to have the correct scaling for spectral convergence of the operator in the limit of large sample size.

We make the following assumptions:

1. The set of  $k$  vectors  $Z = \{z_i\}_{i=1}^N$  was sampled from a manifold in  $R^K$ .
2. Each sample is drawn from an unknown distribution  $\mu(z)$ .
3. The graph Laplacian is a graph representation of the integrating kernel  $w(x, y)$  with vertex set  $V$ .
4. Each vector in  $Z$  is assigned a vertex and weighted edges  $w(x, y)$  between every  $x, y \in Z$ .

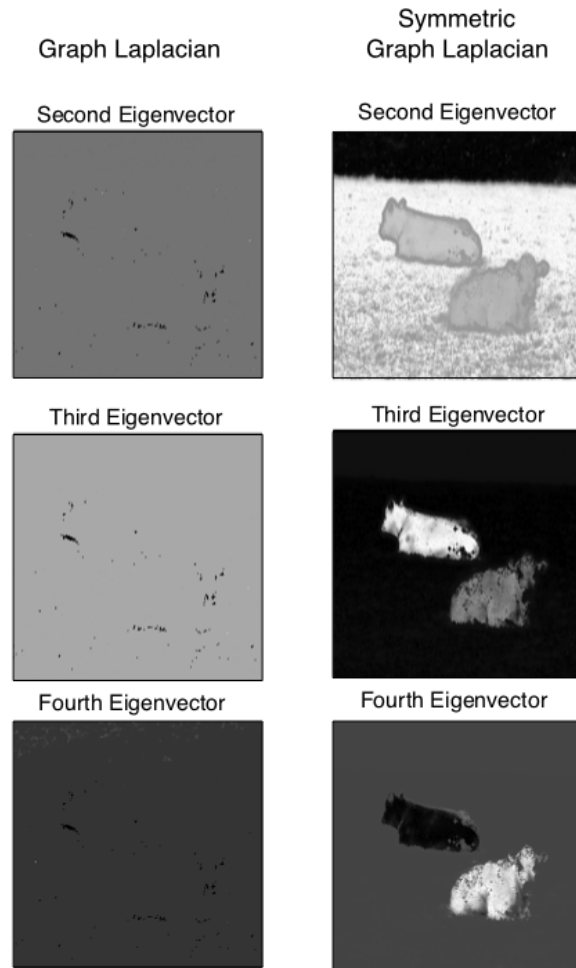
Consistency and practicality of the method require similar and useful solutions as the number of samples increases [85, 83, 84]. Furthermore, the computational methods must be stable. The stability of the computational methods will be discussed first.

Note that the eigenvectors of the discrete Laplacian converge to the eigenvectors of the Laplacian; i.e., the discrete Fourier modes converge to the continuous Fourier modes. Similarly, it has been shown that the spectrum of the graph Laplacian converges (compactly) to the corresponding integral operator [84]. We note that, as stated in [83], there is a dilemma with the convergence for clustering applications. In summary, the unnormalized Laplacian converges to the operator  $L$  defined by  $(Lu)(x) = d(x)u(x) - \int_{\Omega} w(x, y)u(y)dy$ , while the normalized Laplacian converges to  $L_s$  defined by  $(L_s u)(x) = u(x) - \int_{\Omega} (w(x, y)/\sqrt{d(x)d(y)})u(y)dy$ . Both operators are a sum of two operators, a multiplication operator and the operator  $w(x, y)$  or  $w(x, y)/\sqrt{d(x)d(y)}$ . The operators with kernels  $w(x, y)$  and  $w(x, y)/\sqrt{d(x)d(y)}$  are compact and thus have a countable spectrum. The operators  $d(x)$  and the identity operator  $\mathbf{1}$  are multiplication operators, but the operator  $d(x)$  has an a priori unknown value, while the identity operator has an isolated eigenvalue. Note that the spectrum of a multiplication is the essential range of the operator  $d(x)$ ; therefore, by perturbation theory results, the essential spectrum of  $L$  is the essential spectrum of  $d(x)$  [85].

Perturbation theory does not imply anything about the convergence of the eigenvalues inside the essential spectrum of the operator  $L$ . Therefore, we do not know if the function  $L$  is consistent if we increase the number of samples. This problem is avoided if the normalized Laplacian is used instead.

This normalization discussion is not pedantic, and the importance of correct normalization is shown in Figure 2.1. The right column shows example eigenvectors of the symmetric graph Laplacian  $L_s$ . Note that the eigenvectors form reasonable segmentation of the images. For example, the second eigenvector distinguishes between the sky and cows, the third eigenvector separates the cows from the background, and the fourth eigenvector separates the two cows. The left column shows example eigenvalues of the unnormalized Laplacian  $L$ . The spectrum of the unnormalized Laplacian (2.4) is dominated by large spikes at a few pixels. In contrast, the eigenfunctions on the normalized symmetric Laplacian (2.6) provide appealing segmentations of the image.

**2.5. Semisupervised Learning (SSL) on Graphs.** There are numerous approaches to SSL using graph theory, and we mention a few that are related to this work. The work of Coifman, Szlam, and others [18, 76] demonstrates techniques for learning classes using a diffusion framework. Their technique implements the geometric diffu-



**Fig. 2.1** *Eigenfunctions from the graph Laplacian obtained from the cow image in section 4.3. The left column shows eigenvectors of the unnormalized Laplacian  $L$  as in (2.4). The right column shows eigenvectors of the symmetric graph Laplacian  $L_s$  as defined in (2.6).*

sion framework with a random walk probability interpretation. Instead of minimizing an energy functional, they find a time  $s$  when the marginal between known classes is maximized and then classify the rest of the samples using this diffusion time  $s$ . The final segmentation is dominated by the smallest eigenvalues of the random walk graph Laplacian. In contrast, our method is based on an extension of a nonlinear geometric segmentation method applied to general graphs rather than lattices embedded in Euclidean space.

The work of Gilboa and Osher [37, 38] is another closely related technique, inspired in part by earlier work of Buades, Coll, and Morel [12] for denoising. They use the graph Laplacian with an explicit forward time-stepping scheme. The explicit time stepping introduces a stiffness constraint (discussed below) that slows the rate of convergence. Furthermore, their algorithm is stopped at an arbitrary stopping time, while the technique proposed here has an automated stopping criterion.

In [38], a nonlinear nonlocal TV-based method is developed which has remarkable results for texture-based inpainting, although the computational time is not so fast. Our method is a different way of approaching this problem by using the GL functional instead of TV and by taking advantage of fast algorithms for the minimization problem by using the Nyström extension for the graph Laplacian.

Now we show how to use the GL energy on graphs in an SSL application. Assume we have data organized in a graphical structure in which each graph node  $z_i \in Z$  corresponds to a data vector  $x_i$  and the weights between the nodes are constructed using a method such as those described in section 2.2. The goal is to perform a binary segmentation of the data on the graph based on a known subset of nodes (possibly very small) which we denote by  $Z_{data}$ . We denote by  $\lambda$  the characteristic function of  $Z_{data}$ :

$$\lambda(z) = \begin{cases} 1 & \text{if } z \in Z_{data}, \\ 0 & \text{otherwise.} \end{cases}$$

The graph segmentation problem automatically finds a decomposition of the vertices  $Z$  into disjoint sets  $Z_{in} \cup Z_{out}$ . These will be computed by assigning  $\pm 1$  to each of the nodes using a variational procedure of minimizing a GL functional. The known data involves a subset of nodes for which  $+1$  or  $-1$  is already assigned and denoted by  $u_0$  in the variational method. We note that in subsequent works using this method the parameter  $\lambda$  has taken on a constant positive value (sometimes different from one) on the training data and zero elsewhere, allowing the user to put a confidence level on the training data, as is done in analogous variational models in image processing.

The GL functional for SSL is

$$(2.10) \quad E(u) = \frac{\epsilon}{2} \langle u, L_s u \rangle + \frac{1}{4\epsilon} \sum_{z \in Z} (u^2(z) - 1)^2 + \sum_{z \in Z} \frac{\lambda(z)}{2} (u(z) - u_0(z))^2.$$

The fidelity term uses a least-squares fit, allowing for a small amount of misclassification (i.e., noisy data) in the information supplied.

**3. Computational Algorithm.** In this section we go into greater detail regarding the numerical scheme used to find the minimizers of the variational problem. There are two main components to the algorithm—the choice of splitting schemes (1.2) and the computation of the basis functions as eigenfunctions of the graph Laplacian. We cover both below.

**3.1. Convex Splitting Scheme.** Our choice of splitting is motivated by prior work on GL-type functionals for image processing with fidelity [21, 8, 7]. First we review the algorithm as it applies to differential operators in the classical GL regularization. An efficient convex splitting scheme can be derived by writing the GL energy with fidelity as

$$E(u) = E_1(u) - E_2(u)$$

with

$$(3.1) \quad E_1(u) = \frac{\epsilon}{2} \int |\nabla u(x)|^2 dx + \frac{c}{2} \int |u(x)|^2 dx,$$

$$(3.2) \quad E_2(u) = -\frac{1}{4\epsilon} \int (u(x)^2 - 1)^2 dx + \frac{c}{2} \int |u(x)|^2 dx - \int \frac{\lambda(x)}{2} (u(x) - u_0(x))^2 dx.$$

Note that the energy  $E_2$  is not strictly concave, but we can choose the constant  $c$  such that it is concave for  $u$  near and in between the potential wells of  $(u^2 - 1)^2$ . This scheme was chosen so the nonlinear term is in the explicit part of the splitting.

Given the above splitting and since the Fourier transform diagonalizes the Laplace operator, the following numerical scheme solves the Euler–Lagrange equations:

$$\begin{aligned} a_k^{(n)} &= \int e^{ikx} u^{(n)}(x) dx, \\ b_k^{(n)} &= \int e^{ikx} (u^{(n)})^3(x) dx, \\ d_k^{(n)} &= \int e^{ikx} \lambda(x) (u^{(n)}(x) - u_0(x)) dx, \\ \mathcal{D}_k &= 1 + dt (\epsilon k^2 + c), \\ a_k^{(n+1)} &= \mathcal{D}_k^{-1} \left[ \left( 1 + \frac{dt}{\epsilon} + c dt \right) a_k^{(n)} - \frac{dt}{\epsilon} b_k^{(n)} - dt \left( d_k^{(n)} \right) \right]. \end{aligned}$$

Note that the  $H^1$  seminorm is convex and thus appears in the convex part of the energy splitting. The first variation of that yields the Laplace operator, which is a stiff operator to have in an evolution equation. The stiffness results because the eigenvalues of the Laplace operator range from order one negative values to minus infinity. Or in the case of a discrete approximation of the Laplace operator, the eigenvalues range from order one to minus one divided by the square of the smallest length scale of resolution (e.g., the spatial grid size in a finite element or finite difference approximation). By projecting onto the eigenfunctions of the Laplacian, we see that there are many different time scales of decay in the spatial operator and all must be resolved numerically in the case of a forward time-stepping scheme. However, when the Laplace operator is evaluated implicitly, at the new time level, one need not resolve the fastest time scales in the time-stepping scheme.

The same time-stepping scheme can be used if the spectral decomposition of the graph Laplacian is used instead of the Laplacian, and we can use the spectral decomposition for any of the graph Laplacians  $L$ ,  $L_w$ , or  $L_s$ . We used the spectrum of  $L_s$  due to the convergence and scaling issues discussed above. The following is a summary of the method as used in this paper.

Decompose the solution  $u^{(n)}$  at each time step according to the known eigenvectors  $\{\phi_k(x)\}$  of  $L_s$ :

$$u^{(n)}(x) = \sum_k a_k^{(n)} \phi_k(x).$$

Likewise we need to decompose the pointwise cube of  $u$  and the fidelity term,

$$\begin{aligned} [u^{(n)}(x)]^3 &= \sum_k b_k^{(n)} \phi_k(x), \\ \lambda(x) \left( u^{(n)}(x) - u_0(x) \right) &= \sum_k d_k^{(n)} \phi_k(x). \end{aligned}$$

Then the algorithm for the next iteration is given in terms of the coefficients for

$$u^{(n+1)}(x) = \sum_k a_k^{(n+1)} \phi_k(x)$$

in terms of its decomposition using the eigenfunctions of  $L_s$  again as a basis for the solution. Define  $\tilde{\lambda}_k$  to be the eigenvalue associated with the eigenfunction  $\phi_k(x)$ , i.e.,  $L_s \phi_k = \tilde{\lambda}_k \phi_k$ ; then the update equation for  $a_k^{(n)}$  is

$$(3.3) \quad \begin{aligned} \mathcal{D}_k &= 1 + dt (\epsilon \tilde{\lambda}_k + c), \\ a_k^{(n+1)} &= \mathcal{D}_k^{-1} \left[ \left( 1 + \frac{dt}{\epsilon} + c dt \right) a_k^{(n)} - \frac{dt}{\epsilon} b_k^{(n)} - dt \left( d_k^{(n)} \right) \right]. \end{aligned}$$

Convex Splitting for the Graph Laplacian

1. Input  $\leftarrow$  an initial function  $u_0$  and the eigenvalue-eigenvector pairs  $(\tilde{\lambda}_k, \phi_k(x))$  for the graph Laplacian  $L_s$  from (2.6).
2. Set convexity parameter  $c$  and interface scale  $\epsilon$  from (3.2).
3. Set the time step  $dt$ .
4. Initialize  $a_k^{(0)} = \int u(x) \phi_k(x) dx$ .
5. Initialize  $b_k^{(0)} = \int [u_0(x)]^3 \phi_k(x) dx$ .
6. Initialize  $d_k^{(0)} = 0$ .
7. Calculate  $\mathcal{D}_k = 1 + dt (\epsilon \tilde{\lambda}_k + c)$ .
8. For  $n$  less than a set number of iterations  $M$ 
  - (a)  $a_k^{(n+1)} = \mathcal{D}_k^{-1} \left[ \left( 1 + \frac{dt}{\epsilon} + c dt \right) a_k^{(n)} - \frac{dt}{\epsilon} b_k^{(n)} - dt d_k^{(n)} \right]$ ,
  - (b)  $u^{(n+1)}(x) = \sum_k a_k^{(n+1)} \phi_k(x)$ ,
  - (c)  $b_k^{(n+1)} = \int [u^{(n+1)}(x)]^3 \phi_k(x) dx$ ,
  - (d)  $d_k^{(n+1)} = \int \lambda(x) (u^{(n+1)}(x) - u_0(x)) \phi_k(x) dx$ .
9. end for
10. Output  $\leftarrow$  the function  $u^{(M)}(x)$ .

This is a generalization of a classical “pseudospectral” scheme for PDEs in which one goes back and forth between the spectral domain (the coefficients  $a_i^{(n)}$ ) and the graph domain in which we evaluate  $u$  directly at every vertex on the graph. The latter must be done in order to compute the cube  $[u^{(n)}(x)]^3$  and the fidelity term  $\lambda(x) (u^{(n)}(x) - u_0(x))$  which can then be projected back to the spectral domain. Here the convex temporal splitting is very important because it effectively removes the stiffness inherent in the diverse time scales that arise from the range of eigenvalues of the graph Laplacian. Our proposed method is useful only if one has a fast method for determining the eigenfunctions  $\phi_k(x)$  and their corresponding eigenvalues. For the case of fully connected graphs we use the Nyström extension reviewed in the next subsection. A detailed discussion of convergence of this class of methods is presented in [55].

**3.2. Nyström Extension for Fully Connected Graphs.** The spectral decomposition of the matrix  $L_s$  is related to the spectral decomposition

$$D^{-1/2} W D^{-1/2} \phi = \xi \phi$$

through the relationship

$$\begin{aligned} L_s \phi &= (1 - D^{-1/2} W D^{-1/2}) \phi \\ &= (1 - \xi) \phi = \tilde{\lambda} \phi. \end{aligned}$$

Therefore, the convex splitting scheme is efficient if the spectral decomposition of the matrix  $D^{-1/2} W D^{-1/2}$  can be quickly found. The matrix  $W$ , however, is a large

matrix, and it cannot be assumed that the matrix will be sparse. We use the Nyström extension discussed by Fowlkes and coauthors [30, 6, 31] to address this issue.

The Nyström method is a technique for performing matrix completion that has been used in a variety of image processing applications, including spectral clustering [66], kernel principle component analysis [23], and fast Gaussian process calculations. Below we review the Nyström method as used in this paper. Although the method is well known in the graph theory community, we include a summary of the ideas here for the benefit of readers not familiar with these techniques (including the PDE community, which may be interested in extending these ideas to general graph problems).

The Nyström method approximates the eigenvalue equation

$$\int_{\Omega} w(y, x)\phi(x) dx = \gamma\phi(y)$$

using a quadrature rule. Recall that a quadrature rule is a technique for finding  $L$  interpolation weights  $a_j(y)$  and a set of  $L$  interpolation points  $X = \{x_j\}$  such that

$$\sum_{j=1}^L a_j(y)\phi(x_j) = \int_{\Omega} w(y, x)\phi(x) dx + E(y),$$

where  $E(x)$  is the error in the approximation. Our model, however, does not allow us to choose the interpolation points, but rather the interpolation points are randomly sampled from some sample space.

Recall that  $Z = \{z_i\}_{i=1}^N$  is the set of nodes on the graph; it also defines an  $N$ -dimensional vector space with  $W$  as a linear operator on that space. In this work, the Nyström method is used to approximate the eigenvalues of the matrix  $W$  with components  $w(z_i, z_j)$ . A key idea used to produce a fast algorithm is to choose a randomly sampled subset  $X = \{x_i\}_{i=1}^L$  of the entire graph  $Z$  to use as interpolation points, and the interpolation weights are the values of the weight function  $a_j(y) = w(y, x_j)$ .

Partition  $Z$  into two sets  $X$  and  $Y$  with  $Z = X \cup Y$  and  $X \cap Y = \emptyset$ . Furthermore, create the set  $X$  by randomly sampling  $L$  points from  $Z$ . Let  $\phi_k(x)$  be the  $k$ th eigenvector for  $W$ . The Nyström extension approximates the value of  $\phi_k(y_i)$ , up to a scaling factor, using the system of equations

$$(3.4) \quad \sum_{x_j \in X} w(y_i, x_j)\phi_k(x_j) = \gamma\phi_k(y_i).$$

This equation cannot be calculated directly since the eigenvectors  $\phi_k(x)$  are not initially known. This problem is overcome by first approximating the eigenvectors of  $W$  with the eigenvectors of a submatrix of  $W$ . These approximate eigenvalues, however, may not be orthogonal. The approximate eigenvectors will then be orthogonalized, and this final set of eigenvectors will serve as an approximation to the eigenvectors of the complete matrix  $W$ . Note that since only a subset of the matrix  $W$  is initially used, only a subset of the eigenvectors can be approximated using this method.

The notation

$$(3.5) \quad W_{XY} = \begin{bmatrix} w(x_1, y_1) & \dots & w(x_1, y_{N-L}) \\ \vdots & \ddots & \vdots \\ w(x_L, y_1) & \dots & w(x_L, y_{N-L}) \end{bmatrix}$$

will be used in this section. Similarly, define the matrices  $W_{XX}$  and  $W_{YY}$  and the vectors  $\phi_X$  and  $\phi_Y$ . The matrix  $W \in \mathbb{R}^K \times \mathbb{R}^K$  and vectors  $\phi \in \mathbb{R}^K$  can be written as

$$W = \begin{bmatrix} W_{XX} & W_{XY} \\ W_{YX} & W_{YY} \end{bmatrix}$$

and  $\phi = [\phi_X^T \ \phi_Y^T]^T$  with  $\phi^T$  denoting the transpose operation.

The spectral decomposition of  $W_{XX}$  is  $W_{XX} = B_X \Gamma B_X^T$ , where  $B_X$  is the eigenvector matrix of  $W_{XX}$  with each column an eigenvector and  $\Gamma = \text{diag}(\gamma_1, \gamma_2, \dots, \gamma_L)$  are the corresponding eigenvalues. The Nyström extension of (3.4) in matrix form using the interpolation points  $X$  is

$$(3.6) \quad B_Y = W_{YX} B_X \Gamma^{-1}.$$

In short, the  $n$  eigenvectors of  $W$  are approximated by  $B = [B_X^T \ (W_{YX} B_X \Gamma^{-1})^T]^T$ . The associated approximation of  $W = B \Gamma B^T$  is

$$W = \begin{bmatrix} W_{XX} & W_{XY} \\ W_{YX} & W_{YX} W_{XX}^{-1} W_{XY} \end{bmatrix}.$$

From this equation, it can be shown that the large matrix  $W_{YY}$  is approximated by

$$W_{YY} \approx W_{YX} W_{XX}^{-1} W_{XY}.$$

As mentioned in [30], the quality of the approximation to  $W_{YY}$  is given by the norm  $\|W_{YY} - W_{YX} W_{XX}^{-1} W_{XY}\|$ , and this is determined by how well  $W_{YY}$  is spanned by the columns of  $W_{XY}$ .

This decomposition is unsatisfactory since the approximate eigenvectors  $\phi_i(x)$  defined above are not orthogonal. This deficiency can be fixed using the following trick. For arbitrary unitary  $A$  and diagonal matrix  $\Xi$ , if

$$(3.7) \quad \Phi = \begin{bmatrix} W_{XX} \\ W_{YX} \end{bmatrix} (B_X \Gamma^{-1/2} B_X^T) (A \Xi^{-1/2}),$$

the matrix  $W$  can be written as  $W = \Phi \Xi \Phi^T$ . We are therefore free to choose  $A$  unitary such that  $\Phi^T \Phi = \mathbf{1}$ . If such a matrix  $A$  can be found, then the matrix  $W$  will be diagonalized using the unitary matrix  $\Phi$ . Define the operator  $Y = A \Xi^{-1/2}$ ; then the proper choice of  $A$  is given through the relationship

$$\Phi^T \Phi = (Y^T)^{-1} W_{XX} Y^{-1} + (Y^T)^{-1} W_{XX}^{-1/2} W_{XY} W_{YX} W_{XX}^{-1/2} Y^{-1}.$$

If  $\Phi^T \Phi = \mathbf{1}$ , then after multiplying the previous equation on the right by  $\Xi^{1/2} A$  and on the left by the transpose, we have

$$(3.8) \quad A^T \Xi A = W_{XX} + W_{XX}^{-1/2} W_{XY} W_{YX} W_{XX}^{-1/2}.$$

Therefore, if the matrix  $W_{XX} + W_{XX}^{-1/2} W_{XY} W_{YX} W_{XX}^{-1/2}$  is diagonalized, then its spectral decomposition can be used to find an approximate orthogonal decomposition of  $W$  with eigenvectors  $\Phi$  given by (3.7).



The matrix  $W$  must also be normalized in order to use  $L_s$  for segmentation. Normalization of the matrix is a straightforward application of (3.7). In particular, let  $\mathbf{1}_K$  be the  $K$ -dimensional unit vector; then define  $[d_X^T d_Y^T]^T$  as

$$\begin{aligned} \begin{bmatrix} d_X \\ d_Y \end{bmatrix} &= \begin{bmatrix} W_{XX} & W_{XY} \\ W_{YX} & W_{YX}W_{XX}^{-1}W_{XY} \end{bmatrix} \begin{bmatrix} \mathbf{1}_K \\ \mathbf{1}_{N-L} \end{bmatrix} \\ &= \begin{bmatrix} W_{XX} \mathbf{1}_K + W_{XY} \mathbf{1}_{N-L} \\ W_{YX} \mathbf{1}_K + (W_{YX}W_{XX}^{-1}W_{XY}) \mathbf{1}_{N-L} \end{bmatrix}. \end{aligned}$$

Let  $A./B$  denote componentwise division between two matrices  $A$  and  $B$  and  $xy^T$  the outer product of two vectors; then the matrices  $W_{XX}$  and  $W_{XY}$  can be normalized by

$$(3.9) \quad \begin{aligned} \overline{W}_{XX} &= W_{XX} ./ (s_X s_X^T), \\ \overline{W}_{XY} &= W_{XY} ./ (s_X s_Y^T), \end{aligned}$$

where  $s_X = \sqrt{d_X}$  and  $s_Y = \sqrt{d_Y}$ .

The Nyström extension can be summarized by the following pseudocode. If the number of points in  $X$  is much smaller than the number of points in  $Y$ , then the complexity of the algorithm is approximately  $O(N)$ .

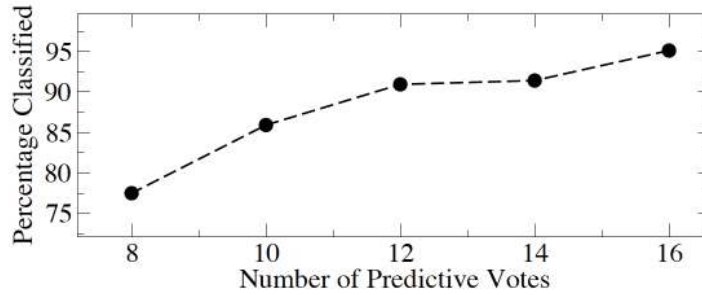
#### Nyström Extension for Symmetric Graph Laplacian

1. Input  $\leftarrow$  a set of features  $Z = \{x_i\}_{i=1}^N$ .
2. Partition the set  $Z$  into  $Z = X \cup Y$ , where  $X$  consists of  $L$  randomly selected elements.
3. Calculate  $W_{XX}$  and  $W_{XY}$  using (3.5).
4.  $d_X = W_{XX} \mathbf{1}_L + W_{XY} \mathbf{1}_{N-L}$ .
5.  $d_Y = W_{YX} \mathbf{1}_L + (W_{YX} W_{XX}^{-1}) (W_{XY} \mathbf{1}_{N-L})$ .
6.  $s_X = \sqrt{d_X}$  and  $s_Y = \sqrt{d_Y}$ .
7.  $\overline{W}_{XX} = W_{XX} ./ (s_X s_X^T)$ .
8.  $\overline{W}_{XY} = W_{XY} ./ (s_X s_Y^T)$ .
9.  $B_X \Gamma B_X^T = \overline{W}_{XX}$  (using the SVD).
10.  $S = B_X \Gamma^{-1/2} B_X^T$ .
11.  $Q = \overline{W}_{XX} + S (\overline{W}_{XY} W_{YX}) S$ .
12.  $A \Xi A^T = Q$  (using the SVD).
13.  $\Phi = \begin{bmatrix} B_X \Gamma^{1/2} \\ W_{YX} B_X \Gamma^{-1/2} \end{bmatrix} B_X^T (A \Xi^{-1/2})$  diagonalizes  $W$ .
14. Output  $\leftarrow$  the  $L$  eigenvalue-eigenvector pairs  $(\phi_i, \tilde{\lambda}_i)$ , where  $\phi_i$  is the  $i$ th column of  $\Phi$  and  $\tilde{\lambda}_i = 1 - \xi_i$  with  $\xi_i$  the  $i$ th diagonal element of  $\Xi$ .

**4. Classification on Graphs.** The GL energy functional can be used for unsupervised and semisupervised classification learning on graphs. This section gives examples of three classification problems. In particular, we investigate the house voting records of 1984 from the UCI machine learning database [32], the Two Moons example dataset of Bühler and Hein [13], and an image segmentation problem.

Each of the classification examples follows the same general procedure. Given a set of vertices  $V = \{\nu_i\}_{i=1}^N$ , the general procedure consists of the following SSL steps:

1. *Determine features:* For each vertex  $\nu_i$ , determine a feature vector  $z_i$ .
2. *Build graph:* Determine edge weights using either formula (2.8) or (2.9) and build an undirected graph based on these weights.



**Fig. 4.1** *The error rate of segmenting the House votes. We tested the accuracy of the segmentation when the most predictive votes were removed. The segmentation procedure was reproduced where we removed the top two, four, six, and eight most predictive votes to investigate the robustness of the algorithm.*

3. *Initialization:* Initialize a function  $u(z_i)$  based on any a priori knowledge.
4. *Minimization:* Minimize the GL energy functional with appropriate constraints and fidelity term(s). Note that for all experiments we use the normalized Laplacian  $L_s$ .
5. *Segmentation:* Segment the vertices into two classes according to  $f(z_i) = \text{sgn}(u(z_i))$ .

Each vertex represents an object in the collection that we want to segment, and the feature vectors provide distinguishing characteristics of the objects.

**4.1. House Voting Records from 1984.** The U.S. House of Representatives voting record dataset consists of 435 individuals, each of whom represents a vertex of the graph. The goal is to use SSL to segment the data by party affiliation, either Democrat or Republican. The SSL algorithm was performed by assuming a party affiliation of five individuals, two Democrats and three Republicans, and segmenting the rest. The votes were taken in 1984 from the 98th United States Congress, 2nd session.

A 16-dimensional feature vector was created using 16 votes recorded for each individual in the following manner. A yes vote was set to one, while a no vote was set to minus one. The voting records had a third category, a did-not-vote category. Each did-not-vote recording was represented by a zero in the feature vector. A fully weighted graph was then created using Gaussian similarity function (2.8) with  $\tau = 0.3$  and with the distance between votes used instead of the square distance.

The function  $u(z)$  was initialized to one for the two Democrats, minus one for the three Republicans, and zero for the rest of the classes. The GL function with fidelity, (2.10), was then minimized using the convex splitting algorithm with parameters  $c = 1$ ,  $dt = 0.1$ ,  $\epsilon = 2$ , and 500 iterations. In the fidelity terms, we chose  $\lambda = 1$  for each of the five known individuals and  $\lambda = 0$  for the rest. This segmentation yielded 95.13% correct results. Note that due to the small size of this graph we did not use the Nyström extension to compute the spectrum.

The probability of the party affiliation given the vote was above 90% for some of the votes. We investigated the accuracy of the segmentation when these votes were removed. Figure 4.1 shows the accuracy of the method when the 14, 12, 10, and 8 least predictive votes were used for the analysis, and we obtained 91.42%, 90.95%, 85.92%, and 77.49% respective accuracy.

The work of Ratanamahatana and Gunopulos [69] studies this dataset using a naive Bayesian decision tree method. They obtain 96.6% classification accuracy using 80% of the data for training and 20% for classification. In contrast, our method uses only 1.15% of the data (5 samples out of 435) to obtain a classification accuracy of 95.13%. The work of Gionis, Mannila, and Tsaparas [39] uses clustering aggregation to automatically determine the number of classes and class membership. Their method obtains 89% correct classification, in contrast to our 95.13%, in which we specify two clusters.

**4.2. Two Moons.** The Two Moons dataset is used by Bühler and Hein [13] and Szlam and Bresson [73, 74] in connection with spectral clustering using the  $p$ -Laplacian. This dataset is constructed using two half circles in  $R^2$  with radius one. The first half circle is contained in the upper half plane with a center at the origin, while the second half circle is created by taking the half circle in the lower half plane and shifting it to  $(1, 0.5)$ . The two half circles are embedded in  $R^{100}$ . Then 2000 data points are sampled from the circles, and independent and identically distributed Gaussian noise with standard deviation 0.02 is added to each of the 100 dimensions. The goal is to segment the two half circles using unsupervised segmentation. The unsupervised segmentation is accomplished by adding a mean zero constraint to the variational problem.

In order to make quantitative comparisons, we build the graph following the procedure described in Szlam and Bresson [73, 74] and Bühler and Hein [13]. They created a 10 nearest neighbor weighted graph from the data using the self-tuning weights of Zelnik-Manor and Perona [90] discussed in section 2.2, where  $M$  was set to 10. This is a difficult segmentation problem since the embedding and noise create a complex graphical structure.

We initialize the function  $u(z)$  using the second eigenfunction of the Laplacian. More specifically, we set  $u(z) = \text{sgn}(\phi_2(z) - \overline{\phi_2})$ , where  $\phi_2(z)$  is the second eigenfunction and  $\overline{\phi_2}$  is the mean of the second eigenfunction. This initialization was used since the second eigenfunction approximates the relaxed graph cut solution (see section 5.1). We minimize the GL energy (2.10) with the mean constraint  $\int u(x)dx = 0$ , but without any fidelity terms. The Nyström extension is ineffective for sparse graphs. Instead, we used the first 20 eigenvectors using the MATLAB sparse matrix algorithms.

Figure 4.2 compares the classical spectral clustering method with our method. Parameters for the GL minimization problem are shown in Table 4.1 and its caption. The left-hand figure is the segmentation achieved by thresholding the eigenvector of the Two Moons dataset. Clearly, spectral clustering using the second eigenvector of the Laplacian does not segment the two half moons accurately. The right-hand image is the segmentation obtained from the algorithm presented in this paper. This algorithm segmented this dataset with an error rate of 2.1%.

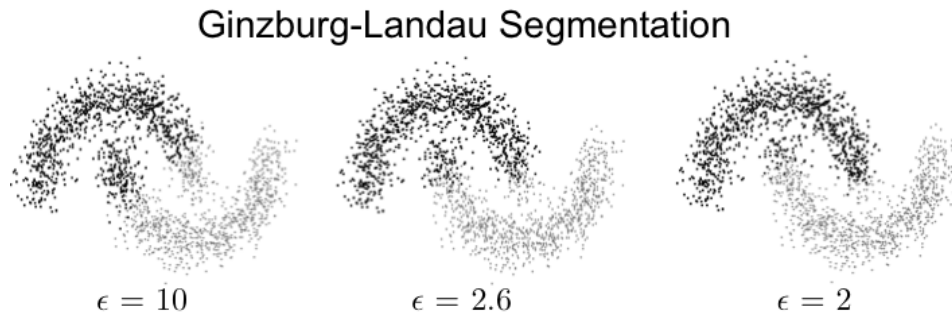
Reducing the GL energy parameter  $\epsilon$  raises the potential barrier between the two states in the GL potential function and reduces the effect of the graph weights. Reducing  $\epsilon$  corresponds to reducing the scale of the graph and allows for a sharper transition between the two states. The change in scale is shown in Figure 4.3, where better segmentation was achieved with reduced  $\epsilon$ . The  $\epsilon = 10$  case is essentially the spectral clustering solution, while the  $\epsilon = 2$  case closely resembles the 1-Laplacian solution of Bühler and Hein [13]. A high quality segmentation in which 94.55% of the samples were classified correctly occurs when the parameter  $\epsilon$  is set to two. This is contrasted with the second eigenvector spectral clustering technique that obtained 82.85% correct classification, essentially equivalent to the large  $\epsilon = 10$  case.



**Fig. 4.2** The left-hand figure is the segmentation achieved by thresholding the second eigenvector of the graph Laplacian. The right-hand image is the segmentation obtained from the algorithm presented in this paper. This algorithm segmented the Two Moons dataset with an error rate of 2.1%.

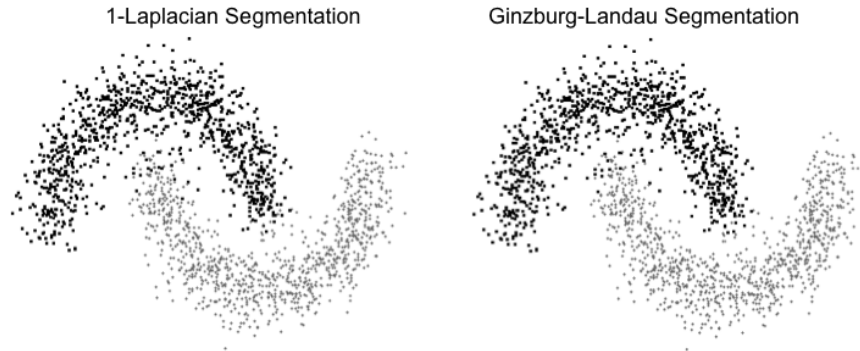
**Table 4.1** Table of parameter values for the GL functional for the Two Moons segmentation. The parameter  $dt = 0.1$  was used throughout along with the formula (2.9) to construct the weighted graph.

$\epsilon$	$c$	No. iterations
10	0.2	500
2	1	200
1.5	1.33	200
1	2	200

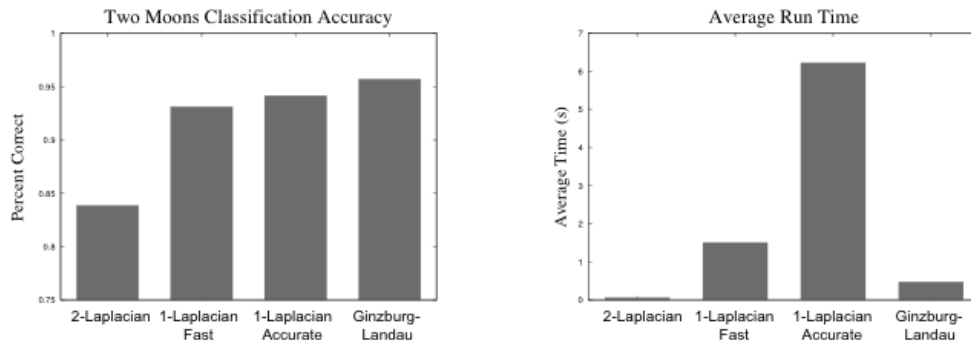


**Fig. 4.3** The parameter  $\epsilon$  determines a scale for the GL energy functional. A more accurate segmentation is obtained as the  $\epsilon$  scale parameter decreases. The percentage of correct classification was 82.85%, 90.75%, and 94.55% for  $\epsilon = 10$ , 2.6, and 2, respectively.

Better segmentation can be achieved if the algorithm is repeated while reducing  $\epsilon$  using the last segmentation as the initialization. The method of successive reductions in  $\epsilon$  was used for image inpainting via the Cahn–Hilliard equation [7, 8]. In [7] the authors carefully study the space of steady states for a stripe inpainting example in which the problem exhibits an incomplete supercritical pitchfork bifurcation as the scale parameter  $\epsilon$  is varied. Different methods of reducing  $\epsilon$  could lead to different local minima of the energy functional, along the two stable branches of the pitchfork. Such a detailed study is beyond the scope of this paper; however, we can examine a segmentation, shown in Figure 4.4, where  $\epsilon$  is reduced from 2 to 1 in steps



**Fig. 4.4** *The 1-Laplacian and the GL clustering methods obtain nearly identical results for unsupervised clustering on the Two Moons dataset. The 1-Laplacian and GL percentage error was 97.3% and 97.9%, respectively.*



**Fig. 4.5** *The accuracy of the GL unsupervised segmentation procedure presented in this paper is compared with spectral clustering and the 1-Laplacian code of Hein and Bühler [42]. These two graphs were created using 100 runs of the Two Moons dataset and averaging the results.*

of 0.5. The final segmentation gives 97.7% accuracy. We compared this segmentation to the 1-Laplacian inverse power method (IPM) of Hein and Bühler [42]. (The code is freely obtainable from [www.ml.uni-saarland.de/code/oneSpectralClustering/oneSpectralClustering.html](http://www.ml.uni-saarland.de/code/oneSpectralClustering/oneSpectralClustering.html).) The normalized 1-Laplacian algorithm of Bühler and Hein with 10 initializations and 5 inner loops was used to obtain 97.3% accuracy for this dataset. The computational time and accuracy of the 1-Laplacian method and the GL technique are shown in Figure 4.5. The GL technique of this paper was able to obtain more accurate results in less computational time. No additional effort was made in our numerical tests to reduce run time—for example, the large number of iterations may not be necessary with an adaptive  $dt$  or a better initialization. Speedup in other problems can easily be an order of magnitude when making such adjustments. Even so, the run time is very fast.

**4.3. Image Labeling.** The objective of image segmentation is to partition an image into multiple components where each component is a set of pixels. Furthermore, each component represents a certain object or class of objects. We are interested in the binary segmentation problem where each pixel can belong to one of two classes.

Most image segmentation algorithms assume that a segmented region is connected in the image. We need not make this assumption. Instead, we build a graph based on feature vectors derived from a neighborhood of each pixel and segment the image based on a partition of the graph. The graph-based segmentation allows us to label the unknown content of one image based on the known content of another image. As input to our segmentation algorithm we take two images, where one of the images has been hand segmented into two classes. The goal is to automatically segment the second image based on the segmentation of the first image.

Each pixel  $y$  represents a vertex of the graph. The feature vector associated with each vertex  $y$  is defined using a pixel neighborhood  $N(y)$  around  $y$ . For example, a typical choice for a pixel neighborhood on a Cartesian grid  $\Omega = Z^2$  is the set

$$N(y) = \{z \in \Omega : |z_1 - y_1| + |z_2 - y_2| \leq M\}$$

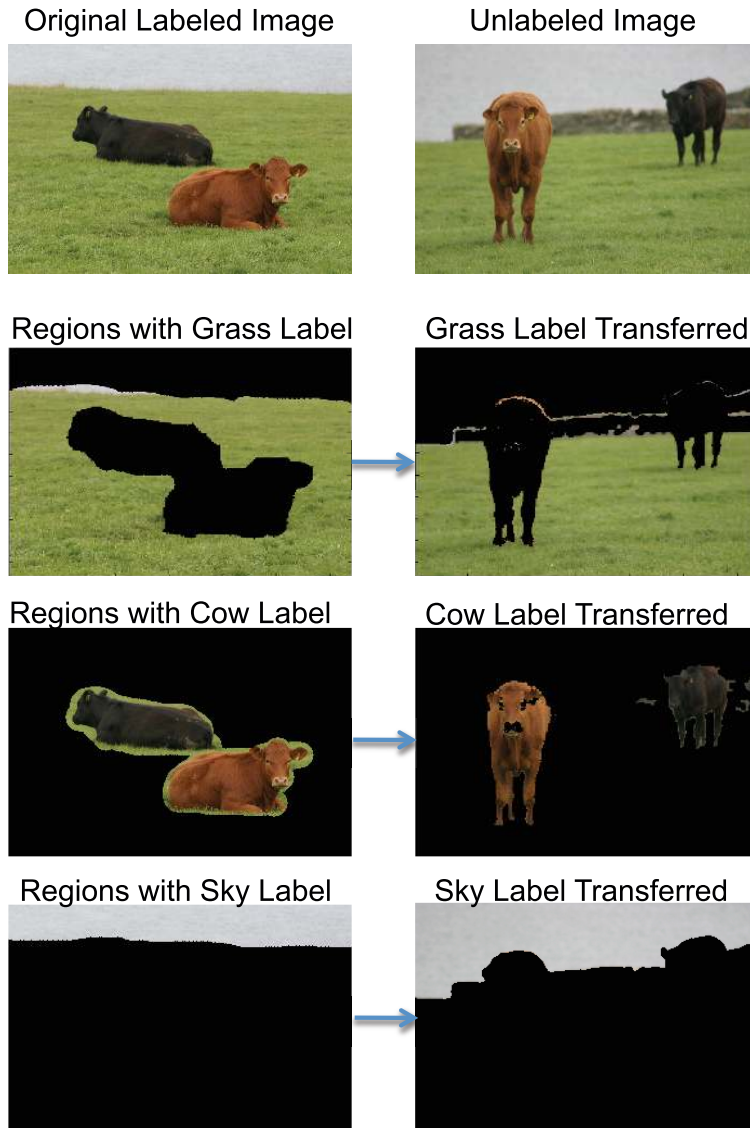
for some integer  $M$ . A feature vector derived from a finite-sized neighborhood of a pixel is called a pixel neighborhood feature.

Let  $I$  be an image; then an example of a pixel neighborhood feature is the set of image pixel values  $z(y) = I(N(y))$  chosen in a consistent order. Another example is to calculate a collection of filter responses for each pixel, i.e.,  $z(x) = ((g_1 * I)(x), (g_2 * I)(x), \dots, (g_j * I)(x))$ , where  $g_i$  represents a filter for each  $i$ , and  $*$  is the convolution operator. The proper choice of neighborhood and features are application dependent. Note that a neighborhood system is equivalent to an edge system in graph theory [15], but the neighborhood system used to determine pixel neighborhood features is *not* the same as the graph used to generate the graph Laplacian.

A fully connected graph is generated using the pixels from two images as vertices and the weight matrix  $w(x, y)$  for edge weights. This graph construction is very general and can be used to segment many different types of objects based on their determining features. For example, color and texture features are appropriate for segmenting trees and grass from other objects. We also note that the metric used in the weight matrix can be modified depending on the dataset. For example, the spectral angle may be more appropriate for segmentation of hyperspectral images.

We demonstrate the image labeling technique on cow images from the Microsoft Research Cambridge Object Recognition Image Database [2] and on face segmentation. The feature vectors used for the Microsoft image database and the face segmentation were, respectively, the Varma–Zisserman MR8 texture features [81] and the van de Weijer–Schmid Hue features [78]. The MR8 texture features are robust to rotation and are translational invariant. The Hue features are invariant to photometric transformations.

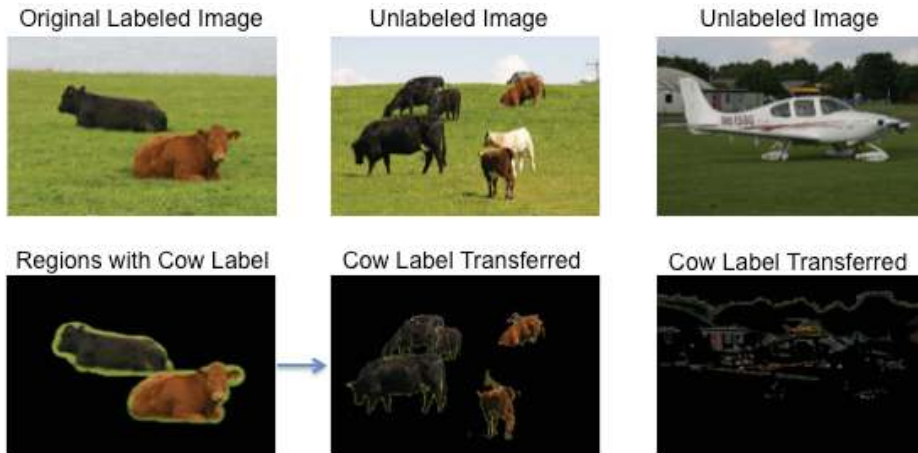
On the hand labeled image, the function  $u(z)$  was initialized to one for one of the classes and minus one for the other class. The function  $u(z)$  was initialized to zero on the unlabeled image. The graph Laplacian was constructed using (2.8). The GL energy with fidelity was then minimized. The parameter values were as follows:  $\tau = 0.1$ ,  $dt = 0.01$ ,  $\epsilon = 0.1$ ,  $c = 21$ , and 500 iterations. The fidelity term  $\Lambda$  was set to one on the initial image and to zero on the unknown image. The Nyström extension was used to determine the spectral decomposition of the weight matrix. The labels of the second image were then determined by the sign of  $u(z)$  on the second image. Results of the segmentation for the Microsoft image database are shown in Figure 4.6. Note that the algorithm is robust to mislabeling in the initial image. Additional examples are shown in Figure 4.7. This figure demonstrates that cows of another color are not labeled as a cow. Note that the white is predominately classified as not a black or



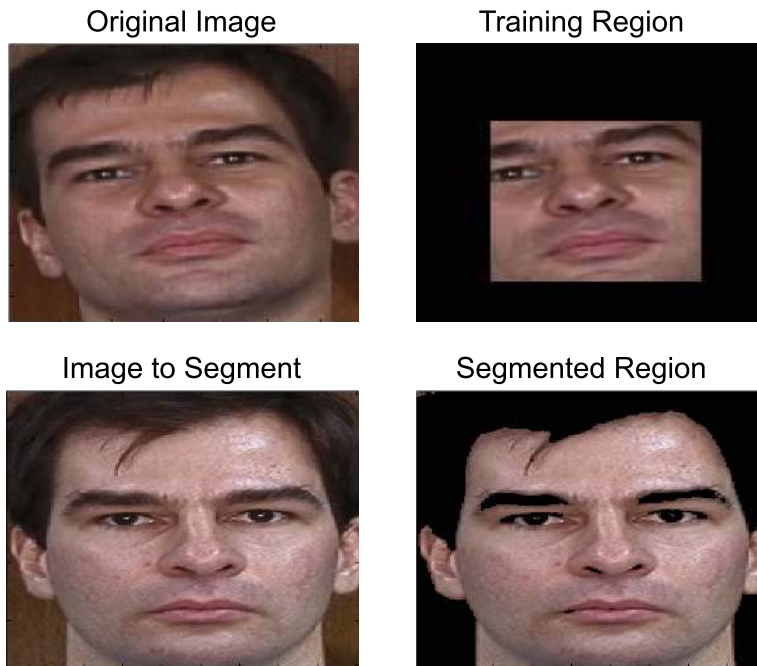
**Fig. 4.6** *The labels from the original upper left-hand image were transferred to the upper right-hand image. The result for each region is shown in the lower images. Note that the algorithm is robust to mislabeled sections. Furthermore, the algorithm can identify regions for which we do not know a label, such as the wall in the right-hand images. Original images in the top row are from the Microsoft Research Cambridge Object Recognition Image Database [2] (copyright © Microsoft Corporation). These images were modified to produce the images in rows 2–4 (8 April 2012).*

brown cow. Furthermore, the algorithm can handle more complex scenes where there are no cows. In the cluttered airplane scene, no consolidated features are identified.

Another example is given in Figure 4.8, where a face from the Georgia Tech Face Database [1] was segmented. In both examples, the predominant features were identified, and some of the pixels with few representative features were removed. For



**Fig. 4.7** This figure uses the same labeled image as in Figure 4.6 and tries to classify similar features in additional images. We include one image with a white colored cow and another image with clutter and no cows. Original images in the top row are from the Microsoft Research Cambridge Object Recognition Image Database [2] (copyright © Microsoft Corporation). These images were modified to produce the images in the bottom row (8 April 2012).



**Fig. 4.8** The robustness of the algorithm to lighting conditions and changes in texture is shown. The upper left-hand image is the original image, and the upper right-hand image is the training region. The lower left-hand image was segmented using the training region in the upper right-hand image, and the segmentation is shown in the lower right-hand image. Original image is from the Georgia Tech Face Database [1] (see also [http://www.anejian.com/research/face\\_reco.htm](http://www.anejian.com/research/face_reco.htm)).



example, the nose and eye of the brown cow were removed from the segmentation, and the eyes and eyebrows of the face were removed. In some other notable image segmentation algorithms a postprocessing filter is often needed. The results shown here use no postprocessing filters.

**5. Connection to Previous Methods in the Literature.** We discuss the connection between our algorithm and related methods from the literature.

**5.1. Graph Cuts.** Spectral clustering and graph segmentation methods are related through the *graph cut* objective function. For two disjoint subsets  $A, B \subset V$ , define the graph cut of two sets as

$$\text{cut}(A, B) = \sum_{x \in A, y \in B} w(x, y) = \sum_x \sum_y |u(x) - u(y)| w(x, y) \equiv |u|_{TV}.$$

The function  $\text{cut}(A, B)$  is smaller when there are fewer weights on the connections between the sets  $A$  and  $B$ . Here  $u$  is the characteristic function that takes value 1 on set  $A$  and 0 on set  $B$ . Thus we see that the size of the graph cut is the same as the total variation of the assignment function  $u$ . This is the equivalence between TV on graphs and the graph cut problem. Hence minimizing the cut is the same as minimizing the total variation of the assignment function. The Gamma convergence result in [79] is the first rigorous connection between graph cuts and the graph GL functional.

The mincut problem involves finding a partition  $\{A, \bar{A}\}$  of  $V$  that minimizes  $\text{cut}(A, \bar{A})$ , where  $\bar{A}$  is the complement of  $A$ . Stoer and Wagner [72] have an efficient algorithm for this problem. The mincut problem, however, leads to poor classification performance for many problems since isolated points often form one cluster and the rest of the points form another cluster. Modifications to the mincut problem include a normalization of either  $|A|$  or  $\text{vol}(A)$  as a measure of the size. This procedure leads to minimization of one of the following [83]:

$$\begin{aligned} \text{RatioCut}(A, \bar{A}) &= \text{cut}(A, \bar{A}) \left( \frac{1}{|A|} + \frac{1}{|\bar{A}|} \right), \\ \text{Ncut}(A, \bar{A}) &= \text{cut}(A, \bar{A}) \left( \frac{1}{\text{vol}(A)} + \frac{1}{\text{vol}(\bar{A})} \right). \end{aligned}$$

**The scalar  $\text{cut}(A, \bar{A})$  is independent of  $|A|$  or  $\text{vol}(A)$**  and the above two sums are smallest when either  $|A|$  or  $\text{vol}(A)$  is the same as  $|\bar{A}|$  or  $\text{vol}(\bar{A})$ , respectively. In other words, the number of vertices or sum of the edge weights must be close to the same in each partition. This balance turns the mincut problem into an NP hard problem [86]. Spectral clustering techniques relax the balancing conditions to approximately solve a simpler version of the mincut problem, and Spielman and Teng have analyzed spectral clustering on bounded-degree planar graphs and finite element meshes [71].

The relaxed minimization of the *RatioCut* is

$$\min_{AC^Y} \langle u, Lu \rangle \text{ such that } u \perp \mathbf{1} \text{ and } \|u\|^2 = |Y|.$$

The relaxed problem is a norm minimization with two constraints, an  $L^2$  constraint and a subspace constraint.

Similar to the *RatioCut* segmentation procedure, the relaxed *NCut* problem is

$$\min_{ACY} \langle u, L_s u \rangle \text{ such that } u \perp D^{1/2} \mathbf{1}, \text{ and } \|u\|^2 = \text{vol}(Y).$$

This minimization problem is in the form of the Rayleigh–Ritz theorem, and the solution is again given by the second eigenvector of  $L_s$ . We emphasize that the difference between the relaxed problem and the true graph cut solution is that the relaxed problem determines a real valued solution, while the graph cut problem finds a binary solution. The relaxation from the discrete problem to the real valued problem does not always yield an approximation to the *NCut* or *RatioCut* problem even for the binary segmentation problem. See, for example, [41] and [77]. The relaxed problem has been used for many segmentation problems, and it produces appealing results [70].

The graph  $p$ -Laplacian is a generalization of the graph Laplacian due to Amghibech [4]. The graph  $p$ -Laplacian is the operator  $L_p$  that satisfies the equation

$$\langle u, L_p u \rangle = \frac{1}{2} \sum_{i,j=1}^N w_{ij} |u_i - u_j|^p.$$

Spectral clustering was accomplished by Bühler and Hein using the graph  $p$ -Laplacian [13]. They defined the eigenvectors of the graph  $p$ -Laplacian using the Rayleigh–Ritz principle where the functional to be minimized is

$$F_p(u) = \frac{\langle u, L_p u \rangle}{\|u\|_p^p}.$$

The work of Szlam and Bresson [73, 74] demonstrated that the solution of the relaxed version of the 1-Laplacian is identical to that of the unrelaxed version. They then derived a split Bregman algorithm to find an approximate solution to the 1-Laplacian. Another approximation method was derived by Bühler and Hein to solve the 1-Laplacian [42]. Their algorithm was used in this work for comparison to the GL segmentation procedure.

The work by Zhu, Ghahramani, and Lafferty [92] is another early publication in this area. Their work connects elements of graph theory, Gaussian random fields, and random walks to provide an outline for supervised learning algorithms using the 2-Laplacian. Their approach constrains  $u$  to be fixed at some of the nodes, but they do not use the double well potential or a fidelity term as is done in this work.

**5.2. Nonlocal Means.** Buades, Coll, and Morel [12] described the following nonlocal filtering, nonlocal means (NLM) procedure for noise removal in images. Define the nonlocal operator

$$NL(u)(x) = \frac{1}{d(x)} \int_{\Omega} w(x, y) u(y) dy,$$

with

$$\begin{aligned} \|u(x) - u(y)\|_a &= \int_{\Omega} G_a(t) |u(x+t) - u(y+t)|^2 dt, \\ (5.1) \quad w(x, y) &= \exp\left(-\frac{\|u(x) - u(y)\|_a}{h^2}\right), \quad d(x) = \int_{\Omega} w(x, y) dy, \end{aligned}$$

and  $G_a(t)$  a Gaussian with standard deviation  $a$ .

Similarly to the segmentation in section 4.3, the norm  $\|\cdot\|_a$  is defined using an image neighborhood. Unlike in section 4.3, the NLM algorithm uses a Gaussian weighted norm so the values of pixels closer to the center pixel have a larger influence on the similarity between two neighborhoods.

Given this analogue with graph theory, the NLM weight matrix can be related to the random walk graph Laplacian. The definition of  $NL(u)$  shows that  $NL(u)(x) = D^{-1}W$ . Substituting  $NL(u)(x) = D^{-1}W$  into (2.7) gives the relationship

$$L_w = \mathbf{1} - NL(u)(x).$$

The operator  $L_w$ , and therefore the NLM operator, naturally occur in the gradient flow of a weighted  $L^2$  norm. To see this, consider the weighted  $L^2$  inner product

$$\langle u, v \rangle_{d(x)} = \int_{\Omega} u(x) v(x) d(x) dx,$$

where  $d(x)$  is the degree function. With this inner product we can write

$$\begin{aligned} \langle u, L_w u \rangle_{d(x)} &= \int_{\Omega} u(x) \left( \int_{\Omega} (u(x) - u(y)) \frac{1}{d(x)} w(x, y) dy \right) d(x) dx \\ &= \int_{\Omega} u(x) \int_{\Omega} (u(x) - u(y)) w(x, y) dy dx \\ &= \langle u, Lu \rangle. \end{aligned}$$

Therefore, there is a natural relationship between the weighted inner product and the nonweighted inner product. This last equation is symmetric when  $x$  and  $y$  are interchanged; therefore we can write the energy functional

$$E(u) = \int_{\Omega} \int_{\Omega} (u(x) - u(y))^2 w(x, y) dy dx = \frac{1}{2} \langle u, (L_w u) \rangle_{d(x)}.$$

Note that the energy  $\langle u, L_w u \rangle_{d(x)}$  is a well-defined energy functional. The gradient flow in the weighted norm  $\langle \cdot \rangle_{d(x)}$  is

$$\frac{\partial u}{\partial t} = -L_w u = -(u(x) - NL(u)(x)).$$

This equation describes a diffusion process using the NLM operator.

Zhou and Scholkopf [91] and Gilboa and Osher [37, 38] derive a calculus based on the nonlocal operators, the former for the discrete graph case and the latter in a continuum setting that was subsequently discretized in computational examples. Zhou and Scholkopf mainly study graph versions of the Poisson equation and its variants. In the continuum case, Gilboa and Osher define the nonlocal derivative for  $y, x \in \Omega$  as

$$(5.2) \quad \partial_y u(x) = (u(y) - u(x)) \sqrt{w(x, y)},$$

where  $0 \leq w(x, y) < \infty$  is the symmetric weight matrix in (2.8). Their nonlocal gradient  $\nabla_w u : \Omega \rightarrow \Omega \times \Omega$  has the form

$$(\nabla_w u)(x, y) = (u(x) - u(y)) \sqrt{w(x, y)}.$$

The nonlocal divergence  $\operatorname{div}_w \vec{v}(x) : \Omega \times \Omega \rightarrow \Omega$  is

$$(\operatorname{div}_w \vec{v})(x) = \int_{\Omega} (v(x, y) - v(y, x)) \sqrt{w(x, y)} dy,$$

which is the adjoint of the nonlocal gradient using the above inner product. Finally, the nonlocal Laplacian can be written as

$$(5.3) \quad \nabla_w^2 u(x) = \frac{1}{2} \operatorname{div}_w (\nabla_w u(x)) = - \int_{\Omega} (u(x) - u(y)) w(x, y) dy.$$

This equation is the continuous equivalent of the standard graph Laplacian, which is a different normalization from the one we use. Gilboa and Osher use this calculus to establish a nonlocal TV energy functional which proves to be highly effective for problems in image inpainting and SSL. In [64] a direct comparison is made between run time for the NLTV model using split Bregman and the MBO scheme on graphs, the results showing similar performance with about a factor of five speedup in run time using MBO. In that work grayscale image inpainting is studied using bitwise binary processing of the gray levels; that step is easily parallelizable. This is in contrast to NLTV, which operates directly on the gray levels as real numbers rather than as bitwise binary signals.

**5.3. Geometric Diffusion.** Coifman et al. [18] discuss a diffusion map formulation to investigate the inherent structure in data and to segment high dimensional datasets. Their construction consists of defining a weight matrix  $w(x, y)$  with admissibility properties satisfied by the Gaussian similarity function  $w(x, y) = \exp(-\|x - y\|^2/\tau)$  used in this paper. A major contribution of geometric diffusion is the observation that a correctly normalized graph Laplacian operator converges to the Laplace–Beltrami operator on a manifold.

A data segmentation procedure was introduced by Coifman et al. using a geometric diffusion approach [18]. The technique was adapted to images by Szlam, Maggioni, and Coifman in [76]. Let  $\Omega_1$  be the set of points in class one, let  $\Omega_2$  be the set of points in class two, and let  $\Omega_3$  be the unlabeled points. Their procedure for a two-class segmentation problem consists of the following steps:

1. Initialize the functions

$$(5.4) \quad u_0^{(i)}(x) = \begin{cases} 1, & x \in \Omega_i, \\ 0 & \text{otherwise.} \end{cases}$$

2. Create the similarity function  $w_{LB}(x, y)$  using feature vectors derived from a neighborhood of each pixel.
3. Diagonalize the matrix  $w_{LB}(x, y) = \sum_j \lambda_j \phi_j(x) \phi_j(y)$ . This step can be performed using the Nystöm extension.
4. Calculate  $u_t^{(i)}(x) = \sum_j \lambda_j^t \phi_j(x) \int \phi_j(y) u_0^{(i)}(y) dy$ , where the parameter  $t$  is chosen by cross-validation with the initial labels.
5. At each point  $x$ , assign the class according to  $\operatorname{argmax}_i \{u_t^{(i)}(x)\}$ .

This equation exploits the result that  $w_{LB}$  is an approximation to the Laplace–Beltrami operator, and therefore  $w_{LB}$  is an approximation to the fundamental solution of the Laplace–Beltrami operator [53]. This method has similarities to an iteration of the MBO scheme on graphs (see [64]) in which the graph assignment function is determined by alternating between a graph heat-equation step and thresholding.

**6. Relevant Recent Work.** This paper developed the first step in the usage of the GL functional for classification of graph-based data. Since its original publication in 2012, this work has inspired many new works, including four completed Ph.D. theses [33, 43, 51, 59] at the time of publication of this SIGEST paper. We summarize the results in published papers here with citations to the relevant references, and discuss research areas where new progress could be made.

**6.1. The MBO Reduction and Multiclass Semisupervised Learning.** Using the original Euclidean GL functional and classical PDE methods, Esedoglu and Tsai developed a simple algorithm for piecewise-constant image segmentation that alternated between evolution of the heat equation and thresholding [28]. That paper built on even earlier work by Merriman, Bence, and Osher (MBO) for motion by mean curvature [57]. Motivated by this work, the MBO computational scheme was extended to the graphical setting by Merkurjev, Kostić, and Bertozzi [64] for binary classification and methods that build on binary classification such as bitwise grayscale classification for inpainting of grayscale images. The graph MBO scheme consists of the following two steps:

1. Heat equation with forcing term. Propagate using

$$\frac{u^{(n+1/2)} - u^{(n)}}{dt} = -L_s u^{(n)} - \lambda(u^{(n)} - u_0).$$

2. Threshold.

$$u^{(n+1)} = \begin{cases} 1 & \text{if } u^{(n+1/2)} \geq 0, \\ 0 & \text{if } u^{(n+1/2)} < 0. \end{cases}$$

The results in [64] showed significant speedup in run time compared to the GL method developed here and also similarly faster run times than the split-Bregman method [40] applied to the Osher–Gilboa nonlocal means graph for the same datasets.

Multiclass versions of both the GL and the graph MBO methods were published in [36, 63] using a symmetric structure for the representation of the classes, building on an earlier multiclass GL method using a modular one-dimensional multiwell functional [34, 35]. The main new idea for multiclass was to rework the binary GL functional to a multiwell functional in  $d$ -dimensions for  $d$  number of classes with the well minima occurring at the corners of the simplex in  $d$  dimensions. The MBO variant involves simply thresholding to the corners of the simplex.

Not only is the extension to multiclass simple, but the algorithm also scales effectively with the number of classes. This scaling is achieved because the bulk of the computational effort is in the eigenfunction calculations of the graph Laplacian which is independent of the class number  $K$ . The iteration scheme depends on  $K$  for the GL and MBO optimization methods, and both methods use the precomputed spectra for the graph heat equation solve. However, the graph heat equation solve can be done independently, and thus in parallel, in each of the  $K$  dimensions. To date, the MBO scheme on graphs has had a bigger impact than the original GL method, although there are still cases in which GL is useful—for example, in a new work on Cheeger cuts [62]. We also note that other approaches using more direct L1 minimization methods for the Cheeger cut problem are developed in [11, 75, 10].

**6.2. Convergence of Graph GL and Related Theory for MBO.** While the theory of Gamma convergence is well known for the classical GL operator and for its wavelet-based cousin [22], until recently, no such theory existed for the graph-based

problem. The first paper to establish such a result was [79] in which the following variant of the GL energy introduced here was proved to Gamma converge to the TV functional on the graph. More precisely we have the following theorem, which is summarized without the infinity cases for simplicity and practicality.

**THEOREM 6.1.** *The GL functional using the original Dirichlet energy given in (2.5) and with the special scaling*

$$\frac{1}{2} \sum_{\mu, \nu \in V} w(\nu, \mu) |u(\nu) - u(\mu)|^2 + \frac{1}{\epsilon} \sum_{\nu \in V} W(u(\nu))$$

*Gamma converges in the epsilon to zero limit to the TV functional on graphs:*

$$\sum_{\mu, \nu \in V} w(\nu, \mu) |u(\nu) - u(\mu)|.$$

We note that this result uses the original graph Laplacian in the Dirichlet inner product rather than the symmetric graph Laplacian  $L_S$  or the random walk version. A detailed connection between these differently scaled Laplacian operators, their associated Dirichlet energies, and their GL Gamma convergence to total variation on graphs are discussed in the two papers [60, 80].

**6.3. Further Comparisons and Global Optimization Methods.** Merkurjev et al. [60] considers the fact that nonconvex functionals such as the GL energy can easily have local minima that do not provide good classification results as compared to the global energy minimizer. On the other hand, one often has to “work” to find machine learning examples using the similarity graph that fail to obtain the global minimizer using the MBO approach for SSL. The examples shown in [60] have very biased sampling of the training data. For example, the samples are from a local quadrant of the MNIST 4-9 dataset when projected onto the plane spanned by two of the lowest eigenfunctions of the graph Laplacian. In contrast, tests with randomly sampled training data gave good results using MBO. That said, it is possible to envision real world applications in which the training data is poorly sampled due to lack of knowledge.

The work in [60] compared performance and run time metrics with the ADMM and max flow L1 convex relaxation methods. The conclusion was that the MBO scheme was notably faster than both; however, the max flow method was notably faster than ADMM. Both of the convex relaxation methods were able to attain the global minimizers for the most difficult test problems in which MBO failed. ADMM was tested with both a spectral truncation method inside the minimization method for the graph Laplace operator versus using a forward method (and thus the full spectrum). The full spectrum method was required to avoid the local minimizers with ADMM. It would be interesting to extend the results of [60] to the multiclass setting.

This paper and the above-mentioned papers all address the SSL problem in which the method is designed to minimize an energy functional of the form of a TV seminorm on the graph plus a least-squares fit to known training data (the semisupervision). We note that in these examples, in contrast to more classical supervised learning methods, the graph cut metric, which as stated above can be approximated by the GL functional or minimized using the MBO scheme, allows for very little training data to achieve good results. For example, in [36] the MNIST dataset is segmented quite accurately with only 3.6% training data compared to the majority of the data used for training for a more classical support vector machine method [52, 20].

**6.4. Modularity Optimization for Unsupervised Clustering.** Unsupervised classification, in which no training data is given and in some problems even the number of classes is unknown, can be a more challenging problem. Spectral clustering is a well-known method for unsupervised classification using graphs in which k-means clustering is applied to the first  $K$  eigenfunctions of the graph Laplacian to provide a segmentation into  $K$  classes. One of most computationally challenging yet accurate methods for unsupervised clustering with unknown numbers of classes is *modularity optimization*, originally developed for community detection on networks (a network is a special case of a weighted undirected graph in which all of the weights are either zero or one) [67]. The modularity function easily extends to weighted undirected graphs, making it an interesting choice for classification of data using similarity graphs. The modularity functional is

$$Q = \frac{1}{m} \sum_{x,y} (w(x,y) - \gamma P(x,y)) \delta(u(x), u(y)),$$

where  $w$  is the standard weight matrix for the graph, and  $u$  is the group assignment function so that  $\delta(u(x), u(y))$  is one if nodes  $x$  and  $y$  are in the same cluster and is zero otherwise.  $P$  is a probability null model; the one used by Newman and Girvan [68] is  $P(x,y) = d(x)d(y)/2m$ , where  $d$  is the degree of the node and  $2m$  is the total volume of the graph, namely,  $\sum_x d(x)$ . Recent work by Hu et al. [44] showed that modularity optimization can be solved using the ideas presented here, giving rise to efficient scalable algorithms for both community detection and unsupervised data classification.

Modularity optimization methods maximize  $Q$  over all possible partitions of the network—a combinatorially intractable problem if one is forced to try all possible partitions. In the case of community detection, the method identifies an optimal clustering without specifying ground truth data or the actual number of classes. Modularity is well studied in the networks literature, but due to its computational complexity, it has not been made useful for machine learning problems on large datasets. This is exemplified by the work [46], in which the GenLouvain algorithm [47] for a multislice variant of modularity optimization [58] was used on the first cow image from this paper. While the code ran for many hours, a very accurate representation of the image was obtained. This work made clear that modularity could be a powerful tool for unsupervised high dimensional data classification, provided faster algorithms could be developed while maintaining accuracy.

For a *fixed number of classes*, the paper [44] demonstrates that maximizing the modularity  $Q$  is equivalent to minimizing a penalized graph cut problem. More precisely, and in terms of total variation, maximizing  $Q$  is equivalent to minimizing

$$(6.1) \quad E_M(u) = |u|_{TV} - \gamma \|u - \bar{u}\|_{\hat{L}^2}^2$$

over all partitions (denoted by the assignment function  $u$ ) with a fixed number of clusters, in which  $\hat{L}^2$  is a degree-weighted  $L^2$  energy (degree of the nodes on the graph) and  $\bar{u}$  is the weighted mean of  $u$ . The constant  $\gamma$  is the same resolution parameter in the original modularity function  $Q$ . In [44] (6.1) was minimized using a variant of the MBO scheme in [64, 36]. This is the first connection between modularity optimization for social networks and L1-compressed sensing methods for data analysis through the graph TV functional. The main advantage to the method in [44] is that it solves precisely the modularity problem constrained to a fixed number of classes, rather than a linear approximation of that problem. When the desired number of classes can be estimated ahead of time, this provides a much faster route to modularity optimization

with sometimes dramatic performance speedup for larger graphs. The method also trivially parallelizes because the optimization for each number of classes can be run simultaneously on different processors, and the partitions subsequently compared.

This idea has potential impact in the analysis of gigantic social network graphs as well as large machine learning problems. For example, as shown in [44], both GenLouvain and the MBO-based methods provide outstanding unsupervised classification of the MNIST dataset.<sup>1</sup> Moreover, the GenLouvain algorithm runs over hours, while the MBO method completes in a couple of minutes to achieve the same results (albeit with searching over a small number of classes, e.g., 2–20 classes to find the optimal 10 or 11 classes, the latter with the one’s divided into digits with and without the flag).

The biggest challenge going forward in adapting our methods to modularity optimization occurs for data that does not have an easily accessible sparse graph realization. The image labeling example presented here is such a case, and we are able to circumvent the burden of computing the  $O(N^2)$  graph for SSL, by performing random sampling using the Nyström extension and using a subset of the eigenfunctions to minimize the GL functional. We note that nowhere in our SSL method are we actually required to compute the graph weights themselves; rather a good approximation of the first variation of the Dirichlet energy is all that is required. Nyström allows us to perform that in approximately  $O(N)$  run time. Likewise a sparse similarity graph would also provide the same computational complexity. In order to apply these ideas to modularity, we need an analogous  $O(N)$  and accurate method of treating the first order differences present in the weighted second moment term in (6.1). This remains a challenge for nonsparse data classification using modularity optimization and is an interesting line of inquiry.

**6.5. Additional Extensions.** Building on the work presented here, simpler unsupervised methods using total variation on graphs have been developed for a fixed number of classes. One example is a multiclass graphical extension of the famous Chan–Vese method for image segmentation [16] applied to hyperspectral video [45]. That work was in part inspired by an earlier paper on SSL using the same large dataset, in which ground truth pixels were unavailable for training data but were approximated by thresholding the first few eigenfunctions of the graph Laplacian [65]. We note that these two papers illustrate that the conventional paradigm of hyperspectral imagery (HSI) analysis, in which a dimension reduction method such as principal component analysis (PCA) is typically done before more complex algorithms are performed, is no longer needed when using these graphical methods. The dimension of the data disappears once the graph is constructed, and the only measurements involving the dimension of the data are done at the level of computing weights of the graph, and in the case of [45], simple averages of the pixel values within each class. In both of these works, an analysis of approximately 280K pixels was carried out on a laptop in a couple of minutes with better results than more standard HSI methods.

An adaptation of the GL method for the Cheeger cut problem is developed in a new work [62]. Another impact of methods outlined in this paper is in the use of auxiliary data for statistical density estimation [87] in which the graph Dirichlet energy associated to the auxiliary data appears as a penalty term for maximum penalized likelihood estimation (MPLE). A specific application studied was crime mapping with high resolution auxiliary data involving residential housing density. The Nyström

---

<sup>1</sup>This calculation is quite surprising to computer scientists more familiar with classical semi-supervised and supervised methods due to the complete lack of training data and the unknown number of classes.



extension was used to compute spectral information about the graph Laplacian which was then used directly in a gradient descent of the MPLE optimization problem. Results were compared with classical kernel density estimation methods and with pure PDE methods without the graphical structure. Yet another work is [24], which takes our original GL approach for SSL and considers the problem of imbalanced classes by playing with the shape of the double well potential. A further extension currently under development is the use of forward solve methods for the heat equation part of MBO such as methods built on the well-known PageRank algorithm [61].

**7. Conclusion.** In summary, this paper develops a class of algorithms for approximating  $L^1$  (TV) regularization for classification of high dimensional data. The algorithms are inspired by classical physical models for diffuse interfaces involving multiple scales, including a diffuse interface scale typically smaller than the bulk features of the problem. Such models have recently been introduced to the image processing literature and have been rigorously connected to methods based on TV. These models are known to produce reasonably sharp edges in image problems, provided that the diffuse interface scale is smaller than the features of interest in the image. By analogy we develop a graph-based method in which the graph Laplacian takes on the role of the spatial Laplace operator in the physical problem. Fast methods can be derived for solving the minimization problem, provided that a reasonably fast algorithm exists for diagonalizing the graph Laplacian. For the classical physics problem there are well-known methods based on the FFT, which diagonalizes the Laplace operator. For our problem we consider standard sparse matrix methods in the case of sparse graphs and the Nystrom extension in the case of highly connected graphs. In all cases we find the results to be comparable to state-of-the-art  $L^1$  methods but with a faster computing time.

In subsequent work there are a few easy steps to speed up the methods described here. One is in regards to sparse graphs: we have used Anderson's Rayleigh–Chebyshev method for computing the eigenfunctions of sparse Hermitian matrices [5]; see [64] for examples using the MBO scheme. Another simple observation is that our iterative method is based on the numerical solution of a gradient descent for a nonlinear functional. We use fixed time steps for this method, and one would expect possibly significant speedup with an efficient variable time step method. One can also exploit variation in the scale parameter  $\epsilon$  during this minimization procedure. This idea was used to great advantage in earlier work on image inpainting using the Cahn–Hilliard equation [8, 7].

**Acknowledgments.** The authors thank Lawrence Carin, John Greer, Gary Hower, Stanley Osher, and Yves van Gennip for comments on the original manuscript and John Urschel for comments on the SIGEST version. Since its original publication, this work has been developed in many directions by a number of students, collaborators and colleagues, including Zach Boyd, Cristina Garcia-Cardona, Yves van Gennip, Tijana Kostić, Huiyi Hu, Thomas Laurent, Kristina Lerman, Xiyang Luo, Zhaoyi Meng, Ekaterina Merkurjev, Allon Percus, Mason Porter, Braxton Osting, Carola-Bibiane Schoenlieb, Justin Sunu, and Joseph Woodworth.

#### REFERENCES

- [1] *Georgia Tech Face Database*, <ftp://ftp.ee.gatech.edu/pub/users/hayes/facedb/>.
- [2] *Microsoft Research Cambridge Object Recognition Image Database, Version 1.0*, <http://research.microsoft.com/downloads> (2005).

- [3] L. AMBROSIO AND V. M. TORTORELLI, *On the approximation of free discontinuity problems*, Boll. Un. Mat. Ital. B (7), 6 (1992), pp. 105–123.
- [4] S. AMGHIBECH, *Eigenvalues of the discrete  $p$ -Laplacian for graphs*, Ars Combin., 67 (2003), pp. 283–302.
- [5] C. ANDERSON, *A Raleigh-Chebyshev procedure for finding the smallest eigenvalues and associated eigenvectors of large sparse Hermitian matrices*, J. Comput. Phys., 229 (2010), pp. 7477–7487.
- [6] S. BELONGIE, C. FOWLKES, F. CHUNG, AND J. MALIK, *Spectral partitioning with indefinite kernels using the Nyström extension*, in Proceedings of the 7th European Conference on Computer Vision (ECCV), Copenhagen, 2002.
- [7] A. BERTOZZI, S. ESEDOĞLU, AND A. GILLETTE, *Analysis of a two-scale Cahn–Hilliard model for binary image inpainting*, Multiscale Model. Simul., 6 (2007), pp. 913–936.
- [8] A. L. BERTOZZI, S. ESEDOĞLU, AND A. GILLETTE, *Inpainting of binary images using the Cahn–Hilliard equation*, IEEE Trans. Image Process., 16 (2007), pp. 285–291.
- [9] A. L. BERTOZZI, N. JU, AND H.-W. LU, *A biharmonic-modified forward time stepping method for fourth order nonlinear diffusion equations*, Discrete Contin. Dyn. Syst., 29 (2011), pp. 1367–1391.
- [10] X. BRESSON, T. LAURENT, D. UMINSKY, AND J. H. VON BRECHT *Convergence and energy landscape for Cheeger cut clustering*, in Advances in Neural Information Processing Systems 25 (NIPS 2012), 2012, pp. 1394–1402.
- [11] X. BRESSON, X.-C. TAI, T. F. CHAN, AND A. SZLAM, *Multi-class transductive learning based on  $\ell^1$  relaxations of Cheeger cut and Mumford-Shah-Potts model*, J. Math. Imag. Vis., 49 (2014), pp. 191–201.
- [12] A. BUADES, B. COLL, AND J. M. MOREL, *A review of image denoising algorithms, with a new one*, Multiscale Model. Simul., 4 (2005), pp. 490–530.
- [13] T. BÜHLER AND M. HEIN, *Spectral clustering based on the graph  $p$ -Laplacian*, in Proceedings of the 26th International Conference on Machine Learning, 2009, pp. 81–88.
- [14] I. CAPUZZO DOLCETTA, S. FINZI VITA, AND R. MARCH, *Area-preserving curve-shortening flows: From phase separation to image processing*, Interfaces Free Bound., 4 (2002), pp. 325–343.
- [15] T. F. CHAN AND J. SHEN, *Image Processing and Analysis: Variational, PDE, Wavelet, and Stochastic Methods*, SIAM, Philadelphia, 2005.
- [16] T. F. CHAN AND L. A. VESE, *Active contours without edges*, IEEE Trans. Image Process., 10 (2001), pp. 266–277.
- [17] F. R. K. CHUNG, *Spectral Graph Theory*, CBMS Reg. Conf. Ser. Math. 92, AMS, Providence, RI, 1997.
- [18] R. R. COIFMAN, S. LAFON, A. B. LEE, M. MAGGIONI, B. NADLER, F. WARNER, AND S. W. ZUCKER, *Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps*, Proc. Natl. Acad. Sci. USA, 102 (2005), pp. 7426–7431.
- [19] J. DARBON AND M. SIGELLE, *A fast and exact algorithm for total variation minimization*, in Pattern Recognition and Image Analysis, Lecture Notes in Comput. Sci. 3522, Springer, Berlin, 2005, pp. 351–359.
- [20] D. DECOSTE AND B. SCHÖLKOPF, *Training invariant support vector machines*, Machine Learning, 46 (2002), pp. 161–190.
- [21] J. A. DOBROSOTSKAYA AND A. L. BERTOZZI, *A wavelet-Laplace variational technique for image deconvolution and inpainting*, IEEE Trans. Image Process., 17 (2008), pp. 657–663.
- [22] J. A. DOBROSOTSKAYA AND A. L. BERTOZZI, *Wavelet analogue of the Ginzburg-Landau energy and its  $\Gamma$ -convergence*, Interfaces Free Bound., 12 (2010), pp. 497–525.
- [23] P. DRINEAS AND M. W. MAHONEY, *On the Nyström method for approximating a Gram matrix for improved kernel-based learning*, J. Mach. Learn. Res., 6 (2005), pp. 2153–2175.
- [24] A. EL GHOUL AND H. SAHBI, *Semi-supervised learning using a graph-based phase field model for imbalanced data set classification*, in the 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2014, pp. 2942–2946, doi:10.1109/ICASSP.2014.6854139.
- [25] S. ESEDOĞLU, *Blind deconvolution of bar code signals*, Inverse Problems, 20 (2004), pp. 121–135.
- [26] S. ESEDOĞLU AND R. MARCH, *Segmentation with depth but without detecting junctions*, J. Math. Imaging Vision, 18 (2003), pp. 7–15.
- [27] S. ESEDOĞLU AND J. SHEN, *Digital inpainting based on the Mumford-Shah-Euler image model*, European J. Appl. Math., 13 (2002), pp. 353–370.
- [28] S. ESEDOĞLU AND Y.-H. R. TSAI, *Threshold dynamics for the piecewise constant Mumford-Shah functional*, J. Comput. Phys., 211 (2006), pp. 367–384.
- [29] D. J. EYRE, *An Unconditionally Stable One-Step Scheme for Gradient Systems*, Technical report, Department of Mathematics, University of Utah, Salt Lake City, UT, 1998.

- [30] C. FOWLKES, S. BELONGIE, F. CHUNG, AND J. MALIK, *Spectral grouping using the Nyström method*, IEEE Trans. Pattern Anal. Mach. Intell., 26 (2004), pp. 214–225.
- [31] C. FOWLKES, S. BELONGIE, AND J. MALIK, *Efficient spatiotemporal grouping using the Nyström method*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2001, pp. 231–238.
- [32] A. FRANK AND A. ASUNCION, *UCI Machine Learning Repository*, School of Information and Computer Science, University of California, Irvine, CA, 2010, <http://archive.ics.uci.edu/ml>.
- [33] C. GARCIA-CARDONA, *Multiclass Learning on Graphs: Diffuse Interface Models and Beyond*, Ph.D. Thesis, Claremont Graduate University, Claremont, CA, 2013.
- [34] C. GARCIA-CARDONA, A. FLENNER, AND A. G. PERCUS, *Multiclass diffuse interface models for semi-supervised learning on graphs*, in Proceedings of the 2nd International Conference on Pattern Recognition Applications and Methods (ICPRAM 2013), 2013, pp. 78–86.
- [35] C. GARCIA-CARDONA, A. FLENNER, AND A. G. PERCUS, *Multiclass semi-supervised learning on graphs using Ginzburg-Landau functional minimization*, Adv. Intell. Syst. Comput., 318 (2015), pp. 119–135.
- [36] C. GARCIA-CARDONA, E. MERKURJEV, A. L. BERTOZZI, A. FLENNER, AND A. G. PERCUS, *Multi-class data segmentation using diffuse interface methods on graphs*, IEEE Trans. Pattern Anal. Mach. Int., 36 (2014), pp. 1600–1613.
- [37] G. GILBOA AND S. OSHER, *Nonlocal linear image regularization and supervised segmentation*, Multiscale Model. Simul., 6 (2007), pp. 595–630.
- [38] G. GILBOA AND S. OSHER, *Nonlocal operators with applications to image processing*, Multiscale Model. Simul., 7 (2008), pp. 1005–1028.
- [39] A. GIONIS, H. MANNILA, AND P. TSAPARAS, *Clustering aggregation*, ACM Trans. Knowl. Discov. Data, 1 (2007), 4.
- [40] T. GOLDSTEIN AND S. OSHER, *The split Bregman method for L1-regularized problems*, SIAM J. Imaging Sci., 2 (2009), pp. 323–343.
- [41] S. GUATTERY AND G. L. MILLER, *On the quality of spectral separators*, SIAM J. Matrix Anal. Appl., 19 (1998), pp. 701–719.
- [42] H. HEIN AND T. BÜHLER, *An inverse power method for nonlinear eigenproblems with applications in 1-spectral clustering and sparse PCA*, in Advances in Neural Information Processing Systems 23 (NIPS 2010), MIT Press, Cambridge, MA, 2010, pp. 847–855.
- [43] H. HU, *Graph Based Models for Unsupervised High Dimensional Data Clustering and Network Analysis*, Ph.D. Thesis, UCLA, 2015.
- [44] H. HU, T. LAURENT, M. A. PORTER, AND A. L. BERTOZZI, *A method based on total variation for network modularity optimization using the MBO scheme*, SIAM J. Appl. Math., 73 (2013), pp. 2224–2246.
- [45] H. HU, J. SUNU, AND A. L. BERTOZZI, *Multi-class graph Mumford-Shah model for plume detection using the MBO scheme*, in Proceedings of the 10th International Conference on Energy Minimization Methods in Computer Vision and Pattern Recognition (Hong Kong), EMMCVPR 2015, X.-C. Tai et al., eds., Lecture Notes in Comput. Sci. 8932, Springer, Berlin, 2015, pp. 209–222.
- [46] H. HU, Y. VAN GENNIP, B. HUNTER, A. L. BERTOZZI, AND M. A. PORTER, *Multislice modularity optimization in community detection and image segmentation*, in IEEE International Conference on Data Mining (ICDM'12), Brussels, 2012, pp. 934–936.
- [47] I. S. JUTLA, L. G. S. JEUB, AND P. J. MUCHA, *A Generalized Louvain Method for Community Detection Implemented in MATLAB*, <http://netwiki.amath.unc.edu/GenLouvain> (2012).
- [48] D. KAY AND R. WELFORD, *A multigrid finite element solver for the Cahn-Hilliard equation*, J. Comput. Phys., 212 (2006), pp. 288–304.
- [49] J. KIM, K. KANG, AND J. LOWENGRUB, *Conservative multigrid methods for Cahn-Hilliard fluids*, J. Comput. Phys., 193 (2004), pp. 551–543.
- [50] R. V. KOHN AND P. STERNBERG, *Local minimisers and singular perturbations*, Proc. Roy. Soc. Edinburgh Sect. A, 111 (1989), pp. 69–84.
- [51] T. KOSTIĆ, *Threshold Dynamics for Statistical Density Estimation and Graph Clustering*, Ph.D. Thesis, UCLA, 2013.
- [52] Y. LECUN, L. BOTTOU, Y. BENGIO, AND P. HAFFNER, *Gradient-based learning applied to document recognition*, Proc. IEEE, 86 (1998), pp. 2278–2324.
- [53] S. S. LAFON, *Diffusion Maps and Geometric Harmonics*, Ph.D. Thesis, Department of Mathematics, Yale University, New Haven, CT, 2004.
- [54] K. LANG, *Fixing two weaknesses of the spectral method*, in Advances in Neural Information Processing Systems 18 (NIPS 2006), MIT Press, Cambridge, MA, 2006, pp. 715–722.
- [55] X. LUO AND A. L. BERTOZZI, *Convergence analysis of the graph Allen-Cahn scheme*, submitted (2016).

- [56] A. MCADAMS, E. SIFAKIS, AND J. TERAN, *A parallel multigrid Poisson solver for fluids simulation on large grids*, in Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA), M. Otaduy and Z. Popovic, eds., 2010, pp. 1–10.
- [57] B. MERRIMAN, J. K. BENCE, AND S. J. OSHER, *Motion of multiple functions: A level set approach*, J. Comput. Phys., 112 (1994), pp. 334–363.
- [58] P. J. MUCHA, T. RICHARDSON, K. MACON, M. A. PORTER, AND J.-P. ONNELA, *Community structure in time-dependent, multiscale, and multiplex networks*, Science, 328 (2010), pp. 876–878, with supplementary material available online, <http://dx.doi.org/10.1126/science.1184819>.
- [59] E. MERKURJEV, *Variational and PDE-Based Methods for Big Data Analysis, Classification and Image Processing Using Graphs*, Ph.D. Thesis, UCLA, 2015.
- [60] E. MERKURJEV, E. BAE, A. L. BERTOZZI, AND X.-C. TAI, *Global binary optimization on graphs for classification of high dimensional data*, J. Math. Imag. Vis., 52 (2015), pp. 414–435.
- [61] E. MERKURJEV, A. L. BERTOZZI, AND F. CHUNG, *A Semi-supervised Heat Kernel Pagerank MBO Algorithm for Classification*, preprint, 2016.
- [62] E. MERKURJEV, A. L. BERTOZZI, K. LERMAN, AND Y. XIAORAN, *Modified Cheeger and ratio cut methods using the Ginzburg-Landau functional for classification of high dimensional data*, submitted (2016).
- [63] E. MERKURJEV, C. GARCIA, A. L. BERTOZZI, A. FLENNER, AND A. PERCUS, *Diffuse interface methods for multiclass segmentation of high-dimensional data*, Appl. Math. Lett., 33 (2014), pp. 29–34.
- [64] E. MERKURJEV, T. KOSTIĆ, AND A. L. BERTOZZI, *An MBO scheme on graphs for classification and image processing*, SIAM J. Imaging Sci., 6 (2013), pp. 1903–1930.
- [65] E. MURKERJEV, J. SUNU, AND A. L. BERTOZZI, *Graph MBO method for multiclass segmentation of hyperspectral stand-off detection video*, in the 2014 IEEE International Conference on Image Processing (ICIP 2014), IEEE, 2014, pp. 689–693.
- [66] M. M. NAEINI, G. DUTTON, K. ROTHLEY, AND G. MORI, *Action recognition of insects using spectral clustering*, in Proceedings of the IAPR Conference on Machine Vision Applications, University of Tokyo, Tokyo, Japan, 2007.
- [67] M. E. J. NEWMAN, *Networks: An Introduction*, Oxford University Press, London, 2010.
- [68] M. E. J. NEWMAN AND M. GIRVAN, *Finding and evaluating community structure in networks*, Phys. Rev. E, 69 (2004), 026113.
- [69] C. A. RATANAMAHATANA AND D. GUNOPULOS, *Scaling up the naive Bayesian classifier: Using decision trees for feature selection*, in Proceedings of the IEEE Workshop on Data Cleaning and Preprocessing (DCAP 2002) at IEEE International Conference on Data Mining (ICDM 2002), Maebashi, Japan, 2002.
- [70] J. SHI AND J. MALIK, *Normalized cuts and image segmentation*, IEEE Trans. Pattern Anal. Mach. Intell., 22 (2000), pp. 888–905.
- [71] D. A. SPIELMAN AND S.-H. TENG, *Spectral partitioning works: Planar graphs and finite element meshes*, Linear Algebra Appl., 421 (2007), pp. 284–305.
- [72] M. STOER AND F. WAGNER, *A simple min-cut algorithm*, J. ACM, 44 (1997), pp. 585–591.
- [73] A. SZLAM AND X. BRESSON, *A Total Variation-Based Graph Clustering Algorithm for Cheeger Ratio Cuts*, UCLA CAM Report 09-68, University of California, Los Angeles, CA, 2009.
- [74] A. SZLAM AND X. BRESSON, *Total variation and Cheeger cuts*, in Proceedings of the 27th International Conference on Machine Learning (ICML-10), J. Fürnkranz and T. Joachims, eds., Omnipress, Haifa, Israel, 2010, pp. 1039–1046.
- [75] A. SZLAM AND X. BRESSON, *Total variation and Cheeger cuts*, in Proceedings of the 27th International Conference on Machine Learning, 2010, pp. 1039–1046.
- [76] A. D. SZLAM, M. MAGGIONI, AND R. R. COIFMAN, *Regularization on graphs with function-adapted diffusion processes*, J. Mach. Learn. Res., 9 (2008), pp. 1711–1739.
- [77] J. C. URSCHER AND T. L. ZIKATANOV, *On the maximal error of spectral approximation of graph bisection*, Linear Multilinear Algebra, 64 (2016), pp. 1–8.
- [78] J. VAN DE WELJER AND C. SCHMID, *Coloring local feature extraction*, in Proceedings of the 9th European Conference on Computer Vision (ECCV 2006), A. Leonardis, H. Bischof, and A. Pinz, eds., Lecture Notes in Comput. Sci. 3952, Springer, Berlin, Heidelberg, 2006, pp. 334–348.
- [79] Y. VAN GENNIP AND A. L. BERTOZZI,  *$\Gamma$ -convergence of graph Ginzburg-Landau functionals*, Adv. Differential Equations, 17 (2012), pp. 1115–1180.
- [80] Y. VAN GENNIP, N. GUILLEN, B. OSTING, AND A. L. BERTOZZI, *Mean curvature, threshold dynamics, and phase field theory on finite graphs*, Milan J. Math., 82 (2014), pp. 3–65.
- [81] M. VARMA AND A. ZISSERMAN, *A statistical approach to texture classification from single images*, Int. J. Comput. Vision, 62 (2005), pp. 61–81.

- [82] P. B. VOLLMAYR-LEE AND A. D. RUTENBERG, *Fast and accurate coarsening simulation with an unconditionally stable time step*, Phys. Rev. E, 68 (2003), 066703.
- [83] U. VON LUXBURG, *A Tutorial on Spectral Clustering*, Technical report TR-149, Max Planck Institute for Biological Cybernetics, Tübingen, Germany, 2006.
- [84] U. VON LUXBURG, M. BELKIN, AND O. BOUSQUET, *On the convergence of spectral clustering on random samples: The normalized case*, in Proceedings of the 17th Annual Conference on Learning Theory (COLT 2004), Banff, Canada, J. Shawe-Taylor and Y. Singer, eds., Lecture Notes in Comput. Sci. 3120, Springer-Verlag, Berlin, Heidelberg, 2004, pp. 457–471.
- [85] U. VON LUXBURG, M. BELKIN, AND O. BOUSQUET, *Consistency of spectral clustering*, Ann. Statist., 36 (2008), pp. 555–586.
- [86] D. WAGNER AND F. WAGNER, *Between min cut and graph bisection*, in Mathematical Foundations of Computer Science, Lecture Notes in Comput. Sci. 711, Springer, Berlin, 1993, pp. 744–750.
- [87] J. T. WOODWORTH, G. O. MOHLER, A. L. BERTOZZI, AND P. J. BRANTINGHAM, *Nonlocal crime density estimation incorporating housing information*, Philos. Trans. A Math. Phys. Eng. Sci., 372 (2014), 20130403.
- [88] L. P. YAROSLAVSKY, *Digital Picture Processing. An Introduction*, Springer-Verlag, Berlin, 1985.
- [89] A. L. YUILLE AND A. RANGARAJAN, *The concave-convex procedure*, Neural Comput., 15 (2003), pp. 915–936.
- [90] L. ZELNIK-MANOR AND P. PERONA, *Self-tuning spectral clustering*, in Advances in Neural Information Processing Systems 17 (NIPS 2004), MIT Press, Cambridge, MA, 2004, pp. 1601–1608.
- [91] D. ZHOU AND B. SCHOLKOPF, *A regularization framework for learning from graph data*, in Proceedings of the ICML Workshop on Statistical Learning and Its Connection to Other Fields, International Conference Proceedings, Vol. 69, ACM, New York, 2004, pp. 132–137.
- [92] X. ZHU, Z. GHAHRAMANI, AND J. LAFFERTY, *Semi-supervised learning using Gaussian fields and harmonic functions*, in Proceedings of the Twentieth International Conference on Machine Learning (ICML), 2003, pp. 912–919.