

A TWO-STAGE CLASSIFICATION METHOD FOR HIGH-DIMENSIONAL DATA AND POINT CLOUDS

XIAOHAO CAI*, RAYMOND CHAN†, XIAOYU XIE‡, AND TIEYONG ZENG‡

Abstract. High-dimensional data classification is a fundamental task in machine learning and imaging science. In this paper, we propose a two-stage multiphase semi-supervised classification method for classifying high-dimensional data and unstructured point clouds. To begin with, a fuzzy classification method such as the standard support vector machine is used to generate a warm initialization. We then apply a two-stage approach named SaT (smoothing and thresholding) to improve the classification. In the first stage, an unconstrained convex variational model is implemented to purify and smooth the initialization, followed by the second stage which is to project the smoothed partition obtained at stage one to a binary partition. These two stages can be repeated, with the latest result as a new initialization, to keep improving the classification quality. We show that the convex model of the smoothing stage has a unique solution and can be solved by a specifically designed primal-dual algorithm whose convergence is guaranteed. We test our method and compare it with the state-of-the-art methods on several benchmark data sets. The experimental results demonstrate clearly that our method is superior in both the classification accuracy and computation speed for high-dimensional data and point clouds.

Key words. Semi-supervised clustering, point cloud classification, variational methods, Graph Laplacian, SaT (smoothing and thresholding).

1. Introduction. Data sets classification is a fundamental task in remote sensing, machine learning, computer vision, and imaging science [1, 56, 43, 52, 38]. The task, simply speaking, is to group the given data into different classes such that, on one hand, data points within the same class shares similar characteristics (e.g. distance, edges, intensities, colors, and textures); on the other hand, pairs of different classes are as dissimilar as possible with respect to certain features. In this paper, we focus on the task of multi-class semi-supervised classification. The total number of classes K of the given data sets is assumed to be known, and a few samples, namely the training points, in each class have been labeled. The goal is therefore to infer the labels of the remaining data points using the knowledge of the labeled ones.

For data classification, previous methods are generally based on graphical models, see e.g. [1, 47, 56], and references therein. In a weighted graph, the data points are vertices and the edge weights signify the affinity or similarity between pairs of data points, where the larger the edge weight is, the closer or more similar the two vertices are. The basic assumption for data classification is that vertices in the graph that are connected by edges with large weight should belong to the same class. Since a fully connected graph is dense and has the size as large as the square of the number of vertices, it is computationally expensive to work on it directly. In order to circumvent this, some cleverly designed approximations have been developed. For example, in [33, 44], spectral approaches are proposed to efficiently calculate the eigendecomposition of a dense graph Laplacian. In [32, 42], the nearest neighbor strategy was adopted to build up a sparse graph where most of its entries are zero, and therefore it is computationally efficient.

*Mullard Space Science Laboratory (MSSL), University College London, Surrey RH5 6NT, UK. Email: x.cai@ucl.ac.uk.

†Department of Mathematics, City University of Hong Kong, Kowloon Tong, Hong Kong. Email: rchan.sci@cityu.edu.hk.

‡Department of Mathematics, The Chinese University of Hong Kong, Shatin, Hong Kong. Emails: xyxie@math.cuhk.edu.hk and zeng@math.cuhk.edu.hk.

In the literature, various studies for semi-supervised classification have been performed by computing a local minimizer of some non-convex energy functional or minimizing a relevant convex relaxation. To name just a few, we have the diffuse interface approaches using phase field representation based on partial differential equation techniques [4, 40], the MBO scheme established for solving the diffusion equation [45, 44, 33], and the variational methods based on graph cut [1, 56]. In particular, in [1], the convex relaxation models and special constraints on class sizes were investigated. In [56], some novelty region-force terms were introduced in the variational models to enforce the affinity between vertices and training samples. To the best of our knowledge, all these proposed variational models have the so-called *no vacuum and overlap constraint* on the labeling functions, which gives rise to non-convex model with NP-hard issues. By allowing labeling functions to take values in the unit simplex, the original NP-hard combinatorial problem is rephrased into a continuous setting, see e.g. [1, 56, 43, 52, 57, 36, 10, 18] for various continuous relaxation techniques (e.g. the ones based on solving the eigenvalue problem, convex approximation, or non-linear optimization) and references therein.

Image segmentation can also be viewed as a special case of the data classification problem [52, 8], since the pixels in an image can be treated as individual points. Various studies and many algorithms have been considered for image segmentation. In particular, variational methods are among the most successful image segmentation techniques, see e.g. [46, 30, 13, 16, 31, 2, 9]. The Mumford-Shah model [46], one of the most important variational segmentation models, was proposed to find piecewise smooth representations of different segments. It is, however, difficult to solve since the model is non-convex and non-smooth. Then substantially rich follow-up works were conducted, and many of them considered compromise techniques such as: (i) simplifying the original complex model, e.g. finding piecewise constant solutions instead of piecewise smooth solutions [27, 41, 55]); (ii) performing convex approximations, e.g. using convex regularization terms like total variation [50, 22]; or (iii) using the smoothing and thresholding (SAT) segmentation methodology [18, 17, 16, 21, 25]; for more details refer to e.g. [26, 35, 11, 39, 26, 48, 49, 58, 5] and references therein. Moreover, various applications were put forward for instance in optical flow [19], tomographic imaging [3], and medical imaging [59, 14, 15, 12, 51, 20].

In this paper, we propose a multi-class semi-supervised data classification method based on the SaT segmentation methodology [18, 17, 16, 21, 25]. It has been shown to be very promising in terms of segmentation quality and computation speed for images corrupted by many different types of blurs and noises. Briefly speaking, the SaT methodology includes two main steps: the first step is to obtain a smooth approximation of the given image through minimizing some convex models; and the second step is to get the segmentation results by thresholding the smooth approximation, e.g. using thresholds determined by the K-means algorithm [37]. Since the models used are convex, the non-convex and NP-hard issues in many existing variational segmentation methods (e.g. the Mumford-Shah model and piecewise constant Mumford-Shah model mentioned above) were naturally avoided.

Our proposed data classification method mainly contains two stages with a warm initialization. The warm initialization is a fuzzy classification result which can be generated by any standard classification methods such as the support vector machine (SVM) [29]; or by labeling the given data randomly if no proper method is available for the given data (e.g. the data set is too large). Its accuracy is not critical since our proposed method will improve the accuracy significantly from this starting point.

With the warm initialization, the first stage of our method is to find a set of smooth labeling functions, where each gives the probability of every point being in a particular class. They are obtained by minimizing a properly-chosen convex objective functional. In detail, the convex objective functional contains K independent convex sub-minimization problems, where each corresponds to one labeling function, with no constraints between these K labeling functions. For each sub-minimization problem, the model is formed by three terms: (i) the data fidelity term restricting the distance between the smooth labeling function and the initialization; (ii) the graph Laplacian (ℓ_2 -norm) term, and (iii) the total variation (ℓ_1 -norm) built on the graph of the given data. The graph Laplacian and the total variation terms regularize the labeling functions to be smooth but at the same time close to a representation on the unit simplex.

After obtaining the set of labeling functions, the second stage of our method is just to project the fuzzy classification results obtained at stage one onto the unit simplex to obtain a binary classification result. This step can be done straightforwardly. To improve the classification accuracy, these two stages can be repeated iteratively, where at each iteration the result at the previous iteration is used as a new initialization.

The main advantage of our proposed method is twofold. First, it performs outstandingly in computation speed, since the proposed model is convex and the K sub-minimization problems are independent with each other (with no constraint on the K labeling functions). The parallelism strategy can be applied straightforwardly to improve computation performance further. On the contrary, the standard start-of-the-art variational data classification methods e.g. [56, 33, 40] have the constraint on unit simplex in their minimization models, so that the non-convex or NP-hard issues can affect seriously the efficiency of these methods, even though some convex relaxations may be applied. Secondly, our method is generally superior in classification accuracy, due to its flexibility of merging the warm initialization and the two-stage iterations which are tractable and manage to improve the accuracy gradually. Note again that we are solving a convex model in the first stage of each iteration, which guarantees a unique global minimizer. In contrast, there is however no guarantee that the results obtained by the standard start-of-the-art variational data classification methods e.g. [56, 33, 40] are global minimizers. The effectiveness of iterations in our proposed method will be shown in the experiments. For most cases, the clustering accuracy would be increased by a significant margin compared to the first initialization and generally outperforms the state-of-the-art variational classification methods.

The paper is organized as follows. In Section 2, we give the basic notation used throughout the paper. In Section 3, we present our method for data sets classification. In Section 4, we present the algorithm for solving the proposed model and its convergence proof. In Section 5, we test our method on benchmark data sets and compare it with the start-of-the-art methods. Conclusions are drawn in Section 6.

2. Basic notation. Let $G = (V, E, w)$ be a weighted undirected graph representing a given point cloud, where V is the vertex set (in which each vertex represents a point) containing N vertices, E is the edge set consisting of pairs of vertices, and $w : E \rightarrow \mathbb{R}_+$ is the weight function defined on the edges in E . The weights $w(\mathbf{x}, \mathbf{y})$ on the edges $(\mathbf{x}, \mathbf{y}) \in E$ measure the similarity between the two vertices \mathbf{x} and \mathbf{y} ; the larger the weight is, the more similar (e.g. closer in distance) the pair of the vertices is.

There are many different ways to define the weight function. Let $d(\cdot, \cdot)$ be a distance metric. Several particularly popular definitions of weight functions are as

follows: (i) radial basis function

$$w(\mathbf{x}, \mathbf{y}) := \exp(-d(\mathbf{x}, \mathbf{y})^2/(2\xi)), \quad \forall(\mathbf{x}, \mathbf{y}) \in E, \quad (2.1)$$

for a prefixed constant $\xi > 0$; (ii) Zelnic-Manor and Perona weight function

$$w(\mathbf{x}, \mathbf{y}) := \exp(-d(\mathbf{x}, \mathbf{y})^2/(\text{var}(\mathbf{x})\text{var}(\mathbf{y}))), \quad \forall(\mathbf{x}, \mathbf{y}) \in E, \quad (2.2)$$

where $\text{var}(\cdot)$ denotes the local variance; and (iii) the cosine similarity

$$w(\mathbf{x}, \mathbf{y}) := \frac{\langle \mathbf{x}, \mathbf{y} \rangle}{\sqrt{\langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle}}, \quad \forall(\mathbf{x}, \mathbf{y}) \in E, \quad (2.3)$$

where $\langle \cdot, \cdot \rangle$ is the inner product.

Let $W = (w(\mathbf{x}, \mathbf{y}))_{(\mathbf{x}, \mathbf{y}) \in E} \in \mathbb{R}^{N \times N}$, the so-called affinity matrix, which is usually assumed to be a symmetric matrix with non-negative entries. Let $D = (h(\mathbf{x}, \mathbf{y}))_{(\mathbf{x}, \mathbf{y}) \in E} \in \mathbb{R}^{N \times N}$ be the diagonal matrix, where its diagonal entries are equal to the sum of the entries on the same row in W , i.e.

$$h(\mathbf{x}, \mathbf{y}) := \begin{cases} \sum_{z \in V} w(\mathbf{x}, \mathbf{z}), & \mathbf{x} = \mathbf{y}, \\ 0, & \text{otherwise.} \end{cases} \quad (2.4)$$

Let $\mathbf{u} = (u(\mathbf{x}))_{\mathbf{x} \in V}^\top \in \mathbb{R}^N$, an N -length column vector. Define the graph Laplacian as $L = D - W$, and the gradient operator ∇ on $u(\mathbf{x}), \forall \mathbf{x} \in V$, as

$$\nabla u(\mathbf{x}) := (w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y})))_{(\mathbf{x}, \mathbf{y}) \in E}. \quad (2.5)$$

Then define the ℓ_1 -norm of an N -length vector as

$$\|\nabla \mathbf{u}\|_1 := \sum_{\mathbf{x} \in V} |\nabla u(\mathbf{x})| = \sum_{(\mathbf{x}, \mathbf{y}) \in E} |w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y}))|, \quad (2.6)$$

and the ℓ_2 -norm (also known as Dirichlet energy)

$$\|\nabla \mathbf{u}\|^2 := \frac{1}{2} \mathbf{u}^\top L \mathbf{u} = \frac{1}{2} \sum_{(\mathbf{x}, \mathbf{y}) \in E} w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y}))^2. \quad (2.7)$$

Note, however, that working with the fully connected graph E —like the setting in (2.5), (2.6) and (2.7)—can be highly computational demanding.

In order to reduce the computational burden, one often only considers the set of edges with large weights. In this paper, the k -nearest-neighbor (k -NN) of a point \mathbf{x} , $\mathcal{N}(\mathbf{x})$, is used to replace the whole edge set starting from the point \mathbf{x} in E . Besides the computational saving, one additional benefit of using k -NN graph is its capability to capture local property of points lying close to a manifold. With the k -NN graph, then the definitions in (2.5), (2.6) and (2.7) become

$$\nabla u(\mathbf{x}) = (w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y})))_{\mathbf{y} \in \mathcal{N}(\mathbf{x})}, \quad (2.8)$$

$$\|\nabla \mathbf{u}\|_1 := \sum_{\mathbf{x} \in V} |\nabla u(\mathbf{x})| = \sum_{\mathbf{x} \in V} \sum_{\mathbf{y} \in \mathcal{N}(\mathbf{x})} |w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y}))|, \quad (2.9)$$

and

$$\|\nabla \mathbf{u}\|^2 := \frac{1}{2} \mathbf{u}^\top L \mathbf{u} = \frac{1}{2} \sum_{\mathbf{x} \in V} \sum_{\mathbf{y} \in \mathcal{N}(\mathbf{x})} w(\mathbf{x}, \mathbf{y})(u(\mathbf{x}) - u(\mathbf{y}))^2, \quad (2.10)$$

respectively, see e.g. [40, 56] for more detail.

3. Proposed data classification method.

3.1. Preliminary. Given a point cloud V containing N points in \mathbb{R}^M . We aim to partition V into K classes V_1, \dots, V_K based on their similarities (the points in the same class possess high similarity), with a set of training points $T = \{T_j\}_{j=1}^K \subset V$, $|T| = N_T$. Note that $T_j \subset V_j$ for $j = 1, \dots, K$. In other words, we aim to assign the points in $V \setminus T$ certain labels between 1 to K using the training set T in which the labels of points are known, and the partition satisfies no vacuum and overlap constraint:

$$V = \bigcup_{j=1}^K V_j \quad \text{and} \quad V_i \cap V_j = \emptyset, \quad \forall i \neq j, 1 \leq i, j \leq K. \quad (3.1)$$

In the rest of the paper, we denote the points in V needed to be labeled as $S = V \setminus T$, and call S the test set in V .

The constraint (3.1) can be described by a binary matrix function $U := (\mathbf{u}_1, \dots, \mathbf{u}_K) \in \mathbb{R}^{N \times K}$ (also called partition matrix), with $\mathbf{u}_j = (u_j(\mathbf{x}))_{\mathbf{x} \in V}^\top \in \mathbb{R}^N : V \rightarrow \{0, 1\}$ defined as

$$u_j(\mathbf{x}) := \begin{cases} 1, & \mathbf{x} \in V_j, \\ 0, & \text{otherwise,} \end{cases} \quad \forall \mathbf{x} \in V, j = 1, \dots, K. \quad (3.2)$$

Clearly, the above definition yields $\sum_{j=1}^K u_j(\mathbf{x}) = 1, \forall \mathbf{x} \in V$. The constraint (3.2) is also known as the indicator constraint on the unit simplex. Since the binary representation in (3.2) generally requires solving a non-convex model with NP-hard issue, a common strategy—the convex unit simplex—is considered as an alternative

$$\sum_{j=1}^K u_j(\mathbf{x}) = 1, \quad \forall \mathbf{x} \in V, \quad \text{s.t.} \quad u_j(\mathbf{x}) \in [0, 1], j = 1, \dots, K. \quad (3.3)$$

Note, importantly, that the convex constraint (3.3) can overcome the NP-hard issue and make some subproblems convex, but generally the whole model can still be non-convex. Therefore, solving a model with constraint (3.1), (3.2), or (3.3) can be time consuming, see e.g. [40, 56] for more detail.

If a result satisfying (3.3) is not completely binary, a common way to obtain an approximate binary solution satisfying (3.2) is to select the binary function as the nearest vertex in the unit simplex by the magnitude of the components, i.e.

$$(u_1(\mathbf{x}), \dots, u_K(\mathbf{x})) \mapsto \mathbf{e}_i, \quad \text{where } i = \underset{j}{\operatorname{argmax}} \{u_j(\mathbf{x})\}_{j=1}^K, \forall \mathbf{x} \in V. \quad (3.4)$$

Here, \mathbf{e}_i is the K -length unit normal vector which is 1 at the i -th component and 0 for all other components.

3.2. Proposed method. In this section, we present our novel two-stage method for data (e.g. point cloud) classification based on the SaT strategy which has been validated very effective in image segmentation. Our method can be summarized briefly as follows: first, a classification result is obtained as a warm initialization by using a classical, fast, but need not be very accurate classification method such as SVM [29]; then, a proposed two-stage iteration scheme is implemented until no change in the labels of the test points could be made between consecutive iterations. Specifically,

at the first stage, we propose to minimize a novel convex model free of constraint (like those in (3.1), (3.2) and (3.3)), to obtain a fuzzy partition, say U , while keeping the training labels unchanged; at the second stage, a binary result is obtained by just applying the binary rule in (3.4) directly on the fuzzy partition obtained at the first stage. This binary result could be the final classification result for the original classification problem, or, if necessary, be set as a new initialization to search a better one in the same manner. In the following, we give the details of each step.

Initialization. Given a point cloud V containing N points in \mathbb{R}^M and training set T containing N_T points with correct labels, we use SVM, which is a standard and fast clustering method as an example, to obtain the first clustering. Let the partition matrix be $\hat{U} = (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_K) \in \mathbb{R}^{N \times K}$, where $\hat{\mathbf{u}}_j = (\hat{u}_j(\mathbf{x}))_{\mathbf{x} \in V}^\top \in \mathbb{R}^N$ for $j = 1, \dots, K$. One could acquire an initialization by any other methods which have better performance than SVM. If no proper method is available (e.g. the data set is too large), then an initialization generated by setting labels to the test points randomly can be used as an alternative.

Stage one. Now we put forward our convex model to find a fuzzy partition U with initialization $\hat{U} = (\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_K)$. It is

$$\operatorname{argmin}_U \sum_{j=1}^K \left\{ \frac{\beta}{2} \|\mathbf{u}_j - \hat{\mathbf{u}}_j\|_2^2 + \frac{\alpha}{2} \mathbf{u}_j^\top L \mathbf{u}_j + \|\nabla \mathbf{u}_j\|_1 \right\}, \quad (3.5)$$

where the first term is the data fidelity term constraining the fuzzy partition not far away from the initialization; the second term is related to $\|\nabla \mathbf{u}\|_2^2$ with graph Laplacian L ; the last term is the total variation constructed on the graph; and $\alpha, \beta > 0$ are regularization parameters. Specifically, the second term in (3.5) is used to impose smooth features on the labels of the points, and the last term is used to force the points with similar information to group together.

We emphasize that we already have the labels on the points in the training set T , with $\bar{U} = (\bar{\mathbf{u}}_1, \dots, \bar{\mathbf{u}}_K) \in \mathbb{R}^{N_T \times K}$ being the partition matrix on T , where $\bar{\mathbf{u}}_j = (\bar{u}_j(\mathbf{x}))_{\mathbf{x} \in T}^\top \in \mathbb{R}^{N_T}$ for $j = 1, \dots, K$. Therefore, we only assign labels to points in the test set S , i.e. we have

$$\hat{u}_j(\mathbf{x}) = \bar{u}_j(\mathbf{x}), \quad \forall \mathbf{x} \in T, j = 1, \dots, K. \quad (3.6)$$

Let $\hat{\mathbf{u}}_{S_j}$ represent the part of $\hat{\mathbf{u}}_j$ defined on the test set S , then we have

$$\hat{\mathbf{u}}_j = (\hat{\mathbf{u}}_{S_j}^\top, \bar{\mathbf{u}}_j^\top)^\top, \quad j = 1, \dots, K. \quad (3.7)$$

We use analogous notations for the partition matrix $U = (\mathbf{u}_1, \dots, \mathbf{u}_K)$, with

$$\mathbf{u}_j = (\mathbf{u}_{S_j}^\top, \bar{\mathbf{u}}_j^\top)^\top, \quad j = 1, \dots, K. \quad (3.8)$$

In Section 4, (3.7) and (3.8) are going to be used to derive an efficient algorithm to solve (3.5).

The following Theorem 3.1 proves that our proposed model (3.5) has a unique solution.

THEOREM 3.1. *Given $\hat{U} \in \mathbb{R}^{N \times K}$ and $\alpha, \beta > 0$, the proposed model (3.5) has a unique solution $U \in \mathbb{R}^{N \times K}$.*

Proof. According to [7, Chapter 9], a strongly convex function has a unique minimum. The conclusion follows easily from the strong convexity of model (3.5). \square

Many algorithms can be used to solve model (3.5) efficiently due to the convexity of the model without constraint. For example, the split-Bregman algorithm [34], which is specifically devised for ℓ_1 regularized problems; the primal-dual algorithm [23], which is designed to solve general saddle point problems; and the powerful alternative, ADMM algorithm [6]. In particular, model (3.5) actually contains K independent sub-minimization problems, where each corresponds to a labeling function \mathbf{u}_j , and therefore the parallelism strategy is ideal to apply. This is an important advantage of our method for large data sets. The algorithm aspects to solve our proposed convex model (3.5) are detailed in Section 4.

Stage two. This stage is to project the fuzzy partition result U obtained at stage one to a binary partition. Here, formula (3.4) is applied to the fuzzy partition U to generate a binary partition, which naturally satisfies no vacuum and overlap constraint (3.1). We remark that compared to the computation time at stage one, the time at stage two is negligible.

Normally, the classification work is complete after we obtain a binary partition matrix at stage two. However, since the way of obtaining an initialization in our scheme is open, and the quality of the initialization could be poor, we suggest going back to stage one with the latest obtained partition as a new initialization and repeat the above two stages until no more change in the partition matrix is observed. More precisely, we set U as \hat{U} and repeat Stages 1 and 2 again to obtain a new U . Then the final classification result is the converged stationary partition matrix, say U^* . Moreover, to accelerate the convergence speed, we update β in (3.5) by a factor of 2 if we are to repeat the stages. This will obviously enforce the closeness between two consecutive clustering results during iterations, which will ensure the algorithm converges fast. We stop the algorithm when no changes are observed in the clustering result compared to the previous one. We remark that a few iterations (≈ 10) are generally enough in practice, see the experimental results in Section 5 for more detail.

Note, importantly, that our classification method here is totally different from other variational methods like [40, 56] which need to minimize variational models with constraint like (3.1), (3.2), (3.3), or other kinds of constraints (e.g. minimum and maximum number of points imposed on individual classes V_i). Even though our proposed model (3.5) has no constraint, the final classification result of our method naturally satisfies no vacuum and overlap constraint (3.1). Therefore, our method is much easier to solve for each iteration. Our proposed method, namely SaT (smoothing and thresholding) method for high-dimensional data classification, is summarized in Algorithm 1.

Algorithm 1 SaT method for high-dimensional data classification

Initialization: Generate initialization \hat{U} by e.g. SVM method.

Output: Binary partition U^* .

For $l = 0, 1, \dots$, until the stopping criterion reached (e.g. $\|U^{(l)} - U^{(l+1)}\| = 0$)

Stage one: Compute fuzzy partition U by solving model (3.5).

Stage two: Compute binary partition $U^{(l+1)}$ by using formula (3.4) on U .

Set $\hat{U} = U^{(l+1)}$ and $\beta = 2\beta$.

Endfor

Set $U^* = U^{(l+1)}$.

4. Algorithm aspects. In this section, we present an algorithm to solve the proposed convex model (3.5) based on the primal-dual algorithm [23] which is briefly

recalled below.

4.1. Primal-dual algorithm. Let X_i be a finite dimensional vector space equipped with a proper inner product $\langle \cdot, \cdot \rangle_{X_i}$ and norm $\|\cdot\|_{X_i}$, $i = 1, 2$. Let map $\mathcal{K} : X_1 \rightarrow X_2$ be a bounded linear operator. The primal-dual algorithm is, generally speaking, to solve the following saddle-point problem

$$\min_{\mathbf{x} \in X_1} \max_{\tilde{\mathbf{x}} \in X_2} \left\{ \langle \mathcal{K}\mathbf{x}, \tilde{\mathbf{x}} \rangle + \mathcal{G}(\mathbf{x}) - \mathcal{F}^*(\tilde{\mathbf{x}}) \right\}, \quad (4.1)$$

where $\mathcal{G} : X_1 \rightarrow [0, +\infty]$, $\mathcal{F} : X_2 \rightarrow [0, +\infty]$ are proper, convex, lower-semicontinuous functions, and \mathcal{F}^* represents the convex conjugate of \mathcal{F} . Given proper initializations, the primal-dual algorithm to solve (4.1) can be summarized in the following iterative way of updating the primal and dual variables

$$\tilde{\mathbf{x}}^{(l+1)} = (I + \sigma \partial \mathcal{F}^*)^{-1}(\tilde{\mathbf{x}}^{(l)} + \sigma \mathcal{K}\mathbf{z}^{(l)}), \quad (4.2)$$

$$\mathbf{x}^{(l+1)} = (I + \tau \partial \mathcal{G})^{-1}(\mathbf{x}^{(l)} - \tau \mathcal{K}^* \tilde{\mathbf{x}}^{(l+1)}), \quad (4.3)$$

$$\mathbf{z}^{(l+1)} = \mathbf{x}^{(l+1)} + \theta(\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}), \quad (4.4)$$

where $\theta \in [0, 1]$, $\tau, \sigma > 0$ are algorithm parameters.

4.2. Algorithm to solve our proposed model. We first define some useful notations which will be used to present our algorithm.

4.2.1. Preliminary. For ease of explanation, in the following, when we say $(i, j) \in E$, the i and j represent the i -th and j -th vertices in E , respectively. Let

$$E' = \{(i, j) \mid i < j, \forall (i, j) \in E\}. \quad (4.5)$$

The graph laplacian $L = D - W \in \mathbb{R}^{N \times N}$ can be decomposed as

$$L = \sum_{(i,j) \in E'} L_{ij}, \quad (4.6)$$

where

$$L_{ij} = \begin{matrix} & & i & & j & & \\ & & \vdots & & \vdots & & \\ i & \left(\begin{array}{ccccc} \cdots & w(i, j) & \cdots & -w(i, j) & \cdots \\ & \vdots & & \vdots & \\ j & \cdots & -w(i, j) & \cdots & w(i, j) & \cdots \\ & \vdots & & \vdots & \end{array} \right) & \in \mathbb{R}^{N \times N} & \\ & & \vdots & & \vdots & & \end{matrix} \quad (4.7)$$

is a matrix with only four nonzero entries which locate at positions (i, i) , (i, j) , (j, i) and (j, j) . Let $E' = E'_a \cup E'_b \cup E'_c$, where

$$E'_a = \{(i, j) \mid i, j \in S, \forall (i, j) \in E'\}, \quad (4.8)$$

$$E'_b = \{(i, j) \mid i, j \in T, \forall (i, j) \in E'\}, \quad (4.9)$$

$$E'_c = E' \setminus (E'_a \cup E'_b). \quad (4.10)$$

Then the decomposition L in (4.6) can be rewritten as

$$L = \sum_{(i,j) \in E'_a} L_{ij} + \sum_{(i,j) \in E'_b} L_{ij} + \sum_{(i,j) \in E'_c} L_{ij}. \quad (4.11)$$

Note that, the terms $\sum_{(i,j) \in E'_a} L_{ij}$ and $\sum_{(i,j) \in E'_b} L_{ij}$ only have nonzero entries which are associated to the test set S and the training set T , respectively. Let

$$\sum_{(i,j) \in E'_a} L_{ij} = \begin{pmatrix} L_S & 0 \\ 0 & 0 \end{pmatrix}, \quad \sum_{(i,j) \in E'_b} L_{ij} = \begin{pmatrix} 0 & 0 \\ 0 & \bar{L} \end{pmatrix}, \quad \sum_{(i,j) \in E'_c} L_{ij} = \begin{pmatrix} L_1 & L_3 \\ L_3^\top & L_2 \end{pmatrix}, \quad (4.12)$$

where $L_S, L_1 \in \mathbb{R}^{(N-N_T) \times (N-N_T)}$ are related to the test set S , $\bar{L}, L_2 \in \mathbb{R}^{N_T \times N_T}$ are related to the training set T , and $L_3 \in \mathbb{R}^{(N-N_T) \times N_T}$. Then we have

$$L = \begin{pmatrix} L_S + L_1 & L_3 \\ L_3^\top & \bar{L} + L_2 \end{pmatrix}. \quad (4.13)$$

According to (2.8), the gradient operator ∇ can be regarded as a linear transformation between \mathbb{R}^N and $\mathbb{R}^{N \times (k-1)}$ (where $k = |\mathcal{N}(\mathbf{x})|$). For a vector $\mathbf{u}_j = (\mathbf{u}_{S_j}^\top, \bar{\mathbf{u}}_j^\top)^\top$ defined in (3.8), let

$$\mathcal{A}_S(\mathbf{u}_{S_j}) = \nabla \begin{pmatrix} \mathbf{u}_{S_j} \\ \mathbf{0} \end{pmatrix} \in \mathbb{R}^{N \times (k-1)}, \quad H_j = \nabla \begin{pmatrix} \mathbf{0} \\ \bar{\mathbf{u}}_j \end{pmatrix} \in \mathbb{R}^{N \times (k-1)}. \quad (4.14)$$

Clearly, $\mathcal{A}_S : \mathbb{R}^{N-N_T} \rightarrow \mathbb{R}^{N \times (k-1)}$ is an operator corresponding to the test set S , and H_j is the gradient matrix corresponding to the training set T which is fixed since $\bar{\mathbf{u}}_j$ is fixed. Then, we have

$$\nabla \mathbf{u}_j = \nabla \begin{pmatrix} \mathbf{u}_{S_j} \\ \bar{\mathbf{u}}_j \end{pmatrix} = \nabla \begin{pmatrix} \mathbf{u}_{S_j} \\ \mathbf{0} \end{pmatrix} + \nabla \begin{pmatrix} \mathbf{0} \\ \bar{\mathbf{u}}_j \end{pmatrix} = \mathcal{A}_S(\mathbf{u}_{S_j}) + H_j. \quad (4.15)$$

4.2.2. Algorithm. Substituting the decomposition of L in (4.13) and ∇ in (4.15), $\hat{\mathbf{u}}_j$ in (3.7) and \mathbf{u}_j in (3.8) into the proposed minimization model (3.5) yields

$$\operatorname{argmin}_{\{\mathbf{u}_{S_j}\}_{j=1}^K} \sum_{j=1}^K \left\{ \frac{\beta}{2} \|\hat{\mathbf{u}}_{S_j} - \mathbf{u}_{S_j}\|_2^2 + \frac{\alpha}{2} \mathbf{u}_{S_j}^\top L_S \mathbf{u}_{S_j} + \alpha \mathbf{u}_{S_j}^\top L_3 \bar{\mathbf{u}}_j + \|\mathcal{A}_S(\mathbf{u}_{S_j}) + H_j\|_1 \right\}. \quad (4.16)$$

Note, obviously, that solving the above model (4.16) is equivalent to solving K sub-minimization problems corresponding to each \mathbf{u}_{S_j} , $j = 1, \dots, K$, which means that our proposed model inherently benefits from the parallelism computation.

For $1 \leq j \leq K$, let

$$\mathcal{G}_j(\mathbf{u}_{S_j}) = \frac{\beta}{2} \|\hat{\mathbf{u}}_{S_j} - \mathbf{u}_{S_j}\|_2^2 + \frac{\alpha}{2} \mathbf{u}_{S_j}^\top L_S \mathbf{u}_{S_j} + \alpha \mathbf{u}_{S_j}^\top L_3 \bar{\mathbf{u}}_j, \quad (4.17)$$

$$\mathcal{F}_j(\tilde{\mathbf{x}}) = \|\tilde{\mathbf{x}} + H_j\|_1. \quad (4.18)$$

Using the definition of the ℓ_1 -norm given in (2.9), the conjugate of \mathcal{F}_j , \mathcal{F}_j^* , can then be calculated as

$$\begin{aligned} \mathcal{F}_j^*(\mathbf{p}) &= \sup_{\tilde{\mathbf{x}} \in \mathbb{R}^{N \times (k-1)}} \langle \tilde{\mathbf{x}}, \mathbf{p} \rangle - \|\tilde{\mathbf{x}} + H_j\|_1 \\ &= -\langle \mathbf{p}, H_j \rangle + \chi_P(\mathbf{p}), \end{aligned} \quad (4.19)$$

where $P = \{\mathbf{p} \in \mathbb{R}^{N \times (k-1)} : \|\mathbf{p}\|_\infty \leq 1\}$, and $\chi_P(\mathbf{p})$ is the characteristic function of set P with value 0 if $\mathbf{p} \in P$, otherwise $+\infty$.

Using the primal-dual formulation in (4.1) with the definitions of \mathcal{G}_j and \mathcal{F}_j^* respectively given in (4.17) and (4.19), then the minimization problem (4.16) corresponding to each \mathbf{u}_{S_j} can be reformulated as

$$\operatorname{argmin}_{\mathbf{u}_{S_j}} \max_{\mathbf{p}} \left\{ \langle \mathcal{A}_S(\mathbf{u}_{S_j}), \mathbf{p} \rangle + \mathcal{G}_j(\mathbf{u}_S) + \langle \mathbf{p}, \mathbf{h}_j \rangle - \chi_P(\mathbf{p}) \right\}. \quad (4.20)$$

To apply the primal-dual method, it remains to compute $(I + \sigma \partial \mathcal{F}_j^*)^{-1}$ and $(I + \tau \partial \mathcal{G}_j)^{-1}$. Firstly, for $\forall \tilde{\mathbf{x}} \in \mathbb{R}^{N \times (k-1)}$, we have

$$\begin{aligned} (I + \sigma \partial \mathcal{F}_j^*)^{-1}(\tilde{\mathbf{x}}) &= \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^{N \times (k-1)}} \mathcal{F}_j^*(\mathbf{p}) + \frac{1}{2\sigma} \|\mathbf{p} - \tilde{\mathbf{x}}\|_2^2 \\ &= \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^{N \times (k-1)}} \chi_P(\mathbf{p}) + \frac{1}{2\sigma} \|\mathbf{p} - \tilde{\mathbf{x}}\|_2^2 - \langle \mathbf{p}, H_j \rangle \\ &= \operatorname{argmin}_{\mathbf{p} \in \mathbb{R}^{N \times (k-1)}} \chi_P(\mathbf{p}) + \frac{1}{2\sigma} \|\mathbf{p} - \tilde{\mathbf{x}} - \sigma H_j\|_2^2 \\ &= \iota_P(\tilde{\mathbf{x}} + \sigma H_j), \end{aligned} \quad (4.21)$$

where the operator $\iota_P(\cdot)$ is the pointwise projection operator onto the set P , i.e., $\forall p \in \mathbb{R}$,

$$\iota_P(p) = \begin{cases} 1, & |p| > 1 \\ p, & \text{otherwise.} \end{cases} \quad (4.22)$$

Secondly, for $\forall \mathbf{x} \in \mathbb{R}^{N-N_T}$, we have

$$(I + \tau \partial \mathcal{G}_j)^{-1}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{u}_{S_j} \in \mathbb{R}^{N-N_T}} \mathcal{G}_j(\mathbf{u}_{S_j}) + \frac{1}{2\tau} \|\mathbf{u}_{S_j} - \mathbf{x}\|_2^2. \quad (4.23)$$

Using the definition of $\mathcal{G}_j(\mathbf{u}_{S_j})$ given in (4.17), problem (4.23) becomes solving the following linear system

$$(\alpha L_S + \beta I + \frac{1}{\tau} I) \mathbf{u}_{S_j} = \beta \hat{\mathbf{u}}_{S_j} + \frac{1}{\tau} \mathbf{x} - \alpha L_3 \bar{\mathbf{u}}_j. \quad (4.24)$$

Since $(\alpha \bar{L} + \beta I + \frac{1}{\tau} I)$ is positive definite, the above linear system can be solved efficiently by e.g. conjugate gradient method [24].

Finally, by exploiting the strong convexity of $\mathcal{G}_j, \forall 1 \leq j \leq K$, which is shown in the next lemma, [23] suggests that we could adaptively modify σ, τ to accelerate the convergence or the primal-dual method.

LEMMA 4.1. *The functions $\mathcal{G}_j, \forall 1 \leq j \leq K$ are strongly convex with parameter β .*

Proof. For simplicity, we omit the subscript j and S_j in the following proof. First, by (4.12), L_S is semi-positive definite. Therefore, $(\frac{\alpha}{2} \mathbf{u}^\top L_S \mathbf{u} + \alpha \mathbf{u}^\top L_3 \bar{\mathbf{u}})$ is convex. Now the strong convexity of \mathcal{G} follows from the fact that the remaining term in (4.17), which is $\frac{\beta}{2} \|\mathbf{u} - \hat{\mathbf{u}}\|_2^2$, is strongly convex with parameter β . \square

The algorithm solving our proposed classification model (4.16) (i.e. model (3.5)) is summarized in Algorithm 2, and its convergence proof is given in Theorem 4.2

below. For each sub-minimization problem, the relative error between two consecutive iterations and/or a given maximum iteration number can be used as stopping criteria to terminate the algorithm. Finally, we emphasize again that our method is quite suitable for parallelism since the K sub-minimization problems are independent with each other and therefore can be computed in parallel.

Algorithm 2 Algorithm solving the proposed model (4.16) (i.e. model (3.5))

Initialization: $\tilde{\mathbf{x}}^{(0)} \in \mathbb{R}^{N \times (k-1)}$, $\mathbf{x}^{(0)}, \mathbf{z}^{(0)} \in \mathbb{R}^{N-N_T}$, $\theta \in [0, 1]$, $\tau^{(0)}, \sigma^{(0)} > 0$.

Output: $\{\mathbf{u}_{S_j}\}_{j=1}^K$.

For $j = 1, \dots, K$ (parallelism strategy can be applied)

For $l = 0, 1, \dots$, until the stopping criterion reached

 Let $\tilde{\mathbf{x}} = \tilde{\mathbf{x}}^{(l)} + \sigma^{(l)} \mathcal{A}_S \mathbf{z}^{(l)}$, and compute $\tilde{\mathbf{x}}^{(l+1)} = (I + \sigma^{(l)} \partial \mathcal{F}^*)^{-1}(\tilde{\mathbf{x}})$ by (4.21);

 Let $\mathbf{x} = \mathbf{x}^{(l)} - \tau^{(l)} \mathcal{A}_S^* \tilde{\mathbf{x}}^{(l+1)}$, and compute $\mathbf{x}^{(l+1)} = (I + \tau^{(l)} \partial \mathcal{G})^{-1}(\mathbf{x})$ by (4.23);

 Let $\theta^{(l)} = 1/\sqrt{1 + \beta \tau^{(l)}}$, and set $\tau^{(l+1)} = \theta^{(l)} \tau^{(l)}$, $\sigma^{(l+1)} = \sigma^{(l)}/\theta^{(l)}$;

 Compute $\mathbf{z}^{(l+1)} = \mathbf{x}^{(l+1)} + \theta(\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)})$;

Endfor

 Set $\mathbf{u}_{S_j} = \mathbf{x}^{(l+1)}$.

Endfor

THEOREM 4.2. *Algorithm 2 converges if $\tau^{(0)}\sigma^{(0)} < \frac{1}{N^2(k-1)}$.*

Proof. By Theorem 2 in [23], Algorithm 2 converges as long as $\|\mathcal{A}_S\|_2^2 < \frac{1}{\tau^{(0)}\sigma^{(0)}}$. Therefore it suffices to find a suitable upper bound for $\|\mathcal{A}_S\|_2$. By our implementation in (4.14) and since the weight functions (eqs. (2.1) to (2.3)) take value between $[-1, 1]$, each entry in \mathcal{A}_S is between $[-1, 1]$. Therefore, the 1-norm and ∞ -norm of \mathcal{A}_S can be easily estimated as

$$\|\mathcal{A}_S\|_1 = \max_{1 \leq j \leq N-N_T} \sum_{i=1}^{N(k-1)} |(\mathcal{A}_S)_{ij}| \leq N(k-1)$$

and

$$\|\mathcal{A}_S\|_\infty = \max_{1 \leq i \leq N(k-1)} \sum_{j=1}^{N-N_T} |(\mathcal{A}_S)_{ij}| \leq N - N_T.$$

Now, we have

$$\|\mathcal{A}_S\|_2 \leq \sqrt{\|\mathcal{A}_S\|_1 \|\mathcal{A}_S\|_\infty} \leq N\sqrt{k-1}.$$

Therefore, we conclude that, the algorithm converges as long as we choose $\tau^{(0)}, \sigma^{(0)} > 0$, such that $\tau^{(0)}\sigma^{(0)} < \frac{1}{N^2(k-1)}$. \square

5. Numerical results. In this section, we evaluate the performance of our proposed method on four benchmark data sets—including THREE MOON, COIL, OPT-DIGITS and MNIST—for semi-supervised learning. THREE MOON is a synthetic data set which has been used frequently e.g. [33, 40, 56]. The COIL, OPT-DIGITS, and MNIST data set can be found from the supplementary material of [28], the UCI machine learning repository¹, and the MNIST Database of Handwritten Digits², respectively. The basic properties of these test data sets are shown in Table 5.1.

¹<http://archive.ics.uci.edu/ml/datasets.html>

²<http://yann.lecun.com/exdb/mnist/>

TABLE 5.1

Basic properties of the test benchmark data sets. “Dimension” means the length of every vector representing individual points in the given data sets.

Data set	Number of classes	Dimension	Number of points
THREE MOON	3	100	1500
COIL	6	241	1500
OPT-DIGITS	10	64	5620
MNIST	10	784	70000

To implement our method, k -NN graphs are constructed for the test data sets, using the randomized kd-tree [53] to find the nearest neighbors with Euclidean distance as the metric. The radial basis function (2.1) is used to compute the weight matrix W , except for the MNIST data set where the Zelnic-Manor and Perona weight function (2.2) is used with eight closed neighbors. The training samples T —samples with labels known—are selected randomly from each test data set. The classification accuracy is defined as the percentage of correctly labeled data points.

For our proposed method, the regularization parameter β is fixed to 10^{-4} for MNIST, 10^{-5} for COIL, and 10^{-2} for THREE MOON, OPT-DIGITS. In practice, one could choose the value of β based on the accuracy of initialization. The better the initialized accuracy, the larger β we could choose. The regularization parameter α is set to 1 for THREE MOON, OPT-DIGITS, 0.4 for MNIST, and 10^{-2} for COIL. The accuracy of the proposed method can be improved further after fine-tuning the values of α and β for individual test data sets. All the codes were run on a MacBook with 2.8 GHz processor and 16 GB RAM, and MATLAB 2017a.

5.1. Methods comparison. As mentioned in previous sections, we use SVM method [29] to generate initializations for our proposed method. If it is not proper for a data set (e.g. very slow due to the large size of the data set), we could just use an initialization generated by assigning clustering labels randomly.

The SVM is a technique aiming to find the best hyperplane that separates data points of one class from the others. In practice, data may not be separable by a hyperplane. In that case, soft margin is used so that the hyperplane would separate many data points if not all. It is also common to kernelize data points, and then find separating hyperplane in the transformed space. The SVM method used in our experiments is trained with linear kernel.

We compare our proposed method with the state-of-the-art methods proposed recently, e.g. CVM [1], GL [33], MBO [33], TVRF [56], LapRF [56], LapRLS [54], MP [54], and SQ-Loss-I [28]. The code TVRF was provided by the authors and the parameters used in it were chosen by trial and error to give the best results. The classification accuracies of methods GL, MBO, LapRF, LapRLS, MP and SQ-Loss-I were taken from [1, 56], in which methods CVM and TVRF were shown to be superior in most of the cases.

5.2. Three moon data. The synthetic THREE MOON data used here is constructed by following the way performed in [1, 56] exactly. We briefly repeat the procedure as follows. First, generate three half circles in \mathbb{R}^2 —two half top unit circles and one half bottom circle with radius 1.5 which are centered at $(0, 0)$, $(3, 0)$ and $(1.5, 0.4)$, respectively. Then 500 points are uniformly sampled from each half circle and embedded into \mathbb{R}^{100} by appending zeros to the remaining dimensions. Finally, an i.i.d Gaussian noise with standard deviation 0.14 is added to each dimension of

the data. An illustration of the first two dimensions of the THREE MOON data is shown in Fig. 5.1 (a) where different colors are applied on each half circle. This is a three-class classification problem with the goal of classifying each half circle using a small number of supervised points from each class. This classification problem is challenging due to the noise and the high dimensionality of all the points with high similarity in \mathbb{R}^{98} .

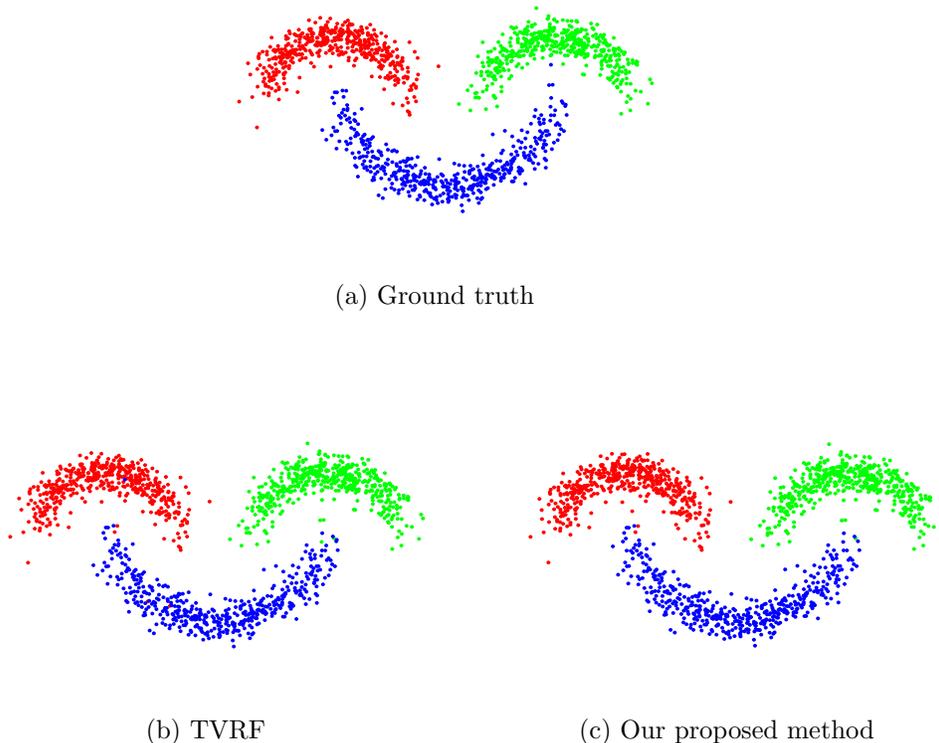


FIG. 5.1. *Three-class classification for THREE MOON synthetic data. (a) Ground truth; (b) and (c) Results of method TVRF [56] and our proposed method, respectively.*

A k -NN graph with $k = 10$ is built for this data set, parameter $\sigma = 3$ is used in the Gaussian weight function, and the distance metric chosen is Euclidean metric for \mathbb{R}^{100} . We first test the methods using uniformly distributed supervised points, where a total number of 75 points are sampled uniformly from this data set as training points.

The accuracies of method TVRF and ours are obtained by running the methods ten times with randomly selected labeled samples, and taking the average of the accuracies. The accuracies of method CVM are copied from the original paper [1]. The accuracy comparison is reported in Table 5.2, which shows that our proposed method gives the highest accuracy; also, see Fig. 5.1 for visual validation of the results between methods of TVRF and ours. The average number of iterations taken for our proposed method is 3.8. Fig. 5.4 (a) gives the convergence history and partition

TABLE 5.2
*Accuracy comparison for THREE
 MOON synthetic data set, with uniformly
 selected training points.*

Method	Accuracy(%)
CVM	98.7
GL	98.4
MBO	99.1
TVRF	98.6
LapRF	98.4
Proposed	99.4

TABLE 5.3
*Accuracy comparison for THREE
 MOON synthetic data set, with non-
 uniformly selected training points.*

Method	Accuracy(%)
TVRF	97.8
Proposed	99.3

accuracy of our proposed method corresponding to iteration steps, which clearly shows the accuracy increment during iterations (note that the accuracy at iteration 0 is the result of the initialization which is obtained by SVM method). Table 5.7 reports the comparison in terms of computation time, which indicates the superior performance of our proposed method in computation speed.

In the following, as a showcase, we test the methods using non-uniformly distributed supervised points, which is used to investigate the robustness of these methods on training points. In this case for the 75 training points, as an example, we respectively pick 5 points from the left and the bottom half circles, and pick the rest 65 points from the right half circle. This sampling is illustrated in Fig 5.2.

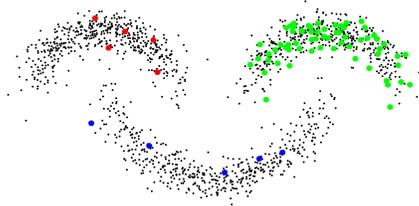


FIG. 5.2. *Unbalanced sampling from THREE MOON data, where sampled points are highlighted with their corresponding labels.*

The accuracies of TVRF and our method are shown in Table 5.3, from which we see clearly that the proposed method gives much higher accuracy. In particular, compared to the results in Table 5.2 using training points chosen uniformly, while the accuracy of TVRF method decreases by 0.8 percent, we observe only a very small decrease in our proposed method. This shows the robustness of our method with respect to the way that training points are selected. Note that in the case of training points chosen non-uniformly, the initialization obtained by SVM is poor, because of which more iterations are needed to converge for our method—average 12.0 iterations in 10 trials versus 3.3 iterations needed for the case of training points selected uniformly.

5.3. COIL data. The benchmark COIL data comes from the Columbia object image library. It contains a set of color images of 100 different objects. These images, with size 128×128 each, are taken from different angles in steps of 5 degrees, i.e., $72 (= 360/5)$ images for each object. In the following, without loss of generality, we also call an image a point for ease of reference. The test data set here is constructed the same way as depicted in e.g. [1, 56] and is briefly described as follows. First, the red channel of each image is down-sampled to 16×16 pixels by averaging over blocks of 8×8 pixels. Then, 24 out of the 100 objects are randomly selected, which amounts to $1728 (= 24 \times 360/5)$ images. After that, these 24 objects are partitioned into six classes with four objects—288 images ($= 4 \times 72$)—in each class. Finally, after discarding 38 images randomly from each class, a data set of 1500 images where 250 images in each of the six classes are constructed. To construct a graph, each image, which is a vector with length 241, is treated as a node on the graph,

For accuracy test, a k -NN graph with $k = 4$ is built for this data set, parameter $\sigma = 250$ is used in the Gaussian weight function, and the distance metric chosen is Euclidean metric for \mathbb{R}^{241} . The training points, amount to 10% of the points, are selected randomly from the data set. Again, we run the test methods 10 times and compare the average accuracy. The resulting accuracies are listed in Table 5.4, showing that our method outperforms other methods. Moreover, the average number of iterations of our method is 12.2. Fig. 5.4 (b) gives the convergence history of our proposed method in partition accuracy corresponding to iterations, which again shows an increasing trend in accuracy.

TABLE 5.4
Accuracy comparison for COIL data set, with uniformly selected training points.

Method	Accuracy(%)
CVM	93.3
TVRF	92.5
LapRF	87.7
GL	91.2
MBO	91.5
Proposed	94.0

TABLE 5.5
Accuracy comparison for MNIST data set, with uniformly selected training points.

Method	Accuracy(%)
CVM	97.7
TVRF	96.9
LapRF	96.9
GL	96.8
MBO	96.9
Proposed	97.5

5.4. MNIST data. The MNIST data set consists of 70,000 images of handwritten digits 0–9, where each image has a size of 28×28 . Fig. 5.3 shows some images of the ten digits from the data set. Each image is a node on a constructed graph. The objective is to classify the data set into 10 disjoint classes corresponding to different digits. For accuracy test, a k -NN graph with $k = 8$ is built for this data set, and Zelnik-Manor and Perona weight function (2.2) is used to compute the weight matrix. The training 2500 (3.57%) points (images) are selected randomly from the total 70,000 points.

The experimental results of the test methods are obtained by running them 10 times with randomly selected training set with a fixed number of points 2500, and the average accuracy is computed for comparison. The accuracies of the test results are shown in Table 5.5, which indicates that our method is comparable to or better than the state-of-the-art methods compared here. Table 5.7 shows the computation time comparison, from which we again see that our method is very competitive in

computation speed. The convergence history of our proposed method in partition accuracy corresponding to iterations is given in Fig. 5.4 (c), which also demonstrates a clear increasing trend in accuracy.

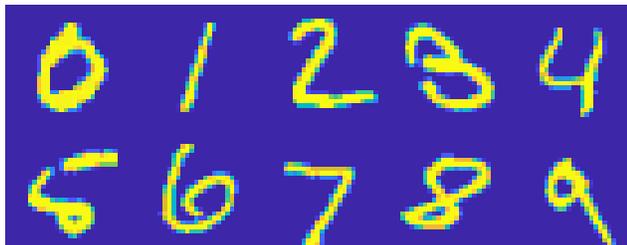


FIG. 5.3. Examples of digits 0–9 from the MNIST data set.

5.5. Opt-Digits data. The OPT-DIGITS data set is constructed as follows. It contains 5620 bitmaps of handwritten digits (i.e. 0–9). Each bitmap has the size of 32×32 and is divided into non-overlapping blocks of 4×4 , and then the number of “on” pixels are counted in each block. Therefore, each bitmap corresponds to a matrix of 8×8 where each element is an integer in $[0, 16]$. The classification problem is to partition the data set into 10 classes.

For accuracy test, a k -NN graph with $k = 8$ is built for this data set, parameter $\sigma = 30$ is used in the Gaussian weight function, and the distance metric chosen is Euclidean metric for \mathbb{R}^{64} . For the experiments on this data set, we generate three training sets respectively with the number of points 50, 100 and 150, which are all selected randomly. All the methods are run 10 times for each training set and the average accuracy is used for comparison. The quantitative results in accuracy are listed in Table 5.6, from which we see that our proposed method is consistently better than the state-of-the-art methods compared for all the cases. We also observe the improvement of the accuracy of these methods w.r.t. the increasing number of points in the training set. Finally, we show the convergence history of our proposed method in partition accuracy corresponding to iterations using the training set with 150 points in Fig. 5.4 (d), which again clearly shows an increasing trend in accuracy.

TABLE 5.6
Accuracy comparison for OPT-DIGITS data set, with uniformly selected training points.

Sample rate	0.89%(50)	1.78%(100)	2.67%(150)
k-NN	85.5	92.0	93.8
SGT	91.4	97.4	97.4
LapRLS	92.3	97.6	97.3
SQ-Loss-I	95.9	97.3	97.7
MP	94.7	97.0	97.1
TVMRF	95.9	98.3	98.2
LapRF	94.1	97.7	98.1
Proposed	96.6	98.5	98.6

5.6. Further discussion. The above experimental results on the benchmark data sets in terms of classification accuracy, shown in Tables 5.2–5.6, indicate that

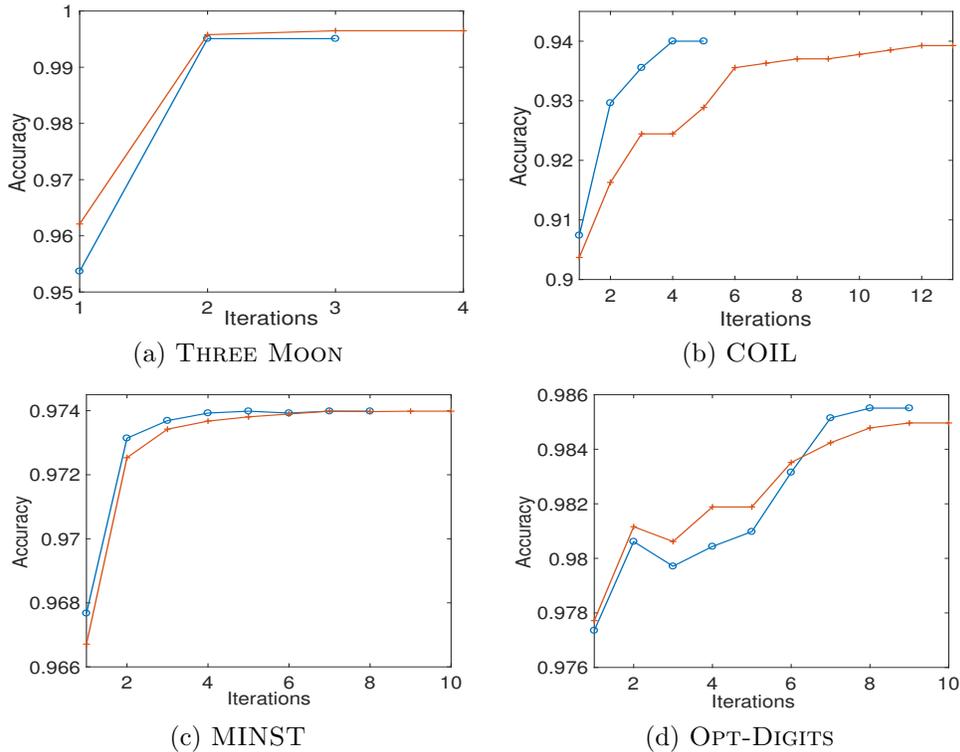


FIG. 5.4. Accuracy convergence history of our proposed method corresponding to iteration steps for all the test data sets; training samples are uniformly selected in each class. Blue curves correspond to cases with the least number of iterations among the 10 trials and orange curves correspond to cases with the largest number of iterations among the 10 trials.

TABLE 5.7

Computation time comparison in seconds. The value in the brackets of our method represents the average number of iterations of the 10 trials. (For Opt-Digits data, we select 100 sample points.)

Method	Computation time in seconds			
	THREE MOON	COIL	MINST	OPT-DIGITS
TVRF	0.71	0.65	66.00	3.42
Proposed	0.30 (3.3)	0.76 (11.7)	82.04 (9.4)	4.45 (9.3)

our proposed method outperforms the state-of-the-art methods for high-dimensional data and point clouds classification.

Compared to the start-of-the-art variational classification models proposed e.g. in [1, 56], in addition to the data fidelity term and ℓ_1 term (e.g. TV), our proposed model in (3.5) contains an additional ℓ_2 term on the labeling functions which is used to smooth the classification results so as to reduce the non-smooth artifact (the so-called staircase artifact in images) introduced by the ℓ_1 term. This is one reason that our method can generally achieve better results. Moreover, the warm initialization used in our method can also play a role to improve the classification quality. Apart from generating the initialization manually, any classification methods can practically be used to generate the initialization. Starting from the initialization, our proposed

method can then be applied to achieve a better classification result by improving the accuracy iteratively. Theoretically speaking, the poorer the quality of the initialization, the more iterations are needed for our method. Nevertheless, we found that even for poor initializations (e.g. the ones generated randomly), 20 iterations are already enough to achieve competitive results. Generally, no more than 15 iterations are needed when using an initialization computed by standard classification methods (e.g. SVM).

Another distinction of our proposed model compared to the variational classification models in e.g. [1, 56] is that there are no constraints on these labeling functions in our objective functional. In other words, in each iteration, we just need to find the minimizer of the objective functional corresponding to each labeling function, but these minimizers do not need to satisfy the constraint that their summation equal to 1. Therefore, the computation speed for every single iteration is improved in our method compared to other methods which have constraints. We emphasize again that, since minimizing each sub-problem with respect to each labeling function is irrelevant to minimizing the sub-problems with respect to other labeling functions, parallelism techniques can be used straightforwardly to further improve the computation performance of our algorithm; theoretically, we just require $1/K$ of the computation time needed for the non-parallelism scheme. This will be extremely important for large data sets. From Table 5.7, we see that, for all the computation time of our method, when considering the effect of parallel computing, our method should be able to outperform the state-of-the-art methods by a large margin.

6. Conclusions. In this paper, a two-stage multiphase semi-supervised method is proposed for classifying high-dimensional data or unstructured point clouds. The method is based on the SaT strategy which has been shown very effective for segmentation problems such as gray or color images corrupted by different degradations. Starting with a proper initialization which can be obtained by using any standard classification algorithm (e.g. SVM) or constructed by users, the first stage of our method is to solve a convex variational model without constraint. Most importantly, our proposed model is a lot easier to solve than the state-of-the-art variational models proposed recently (e.g. [1, 56]) for point clouds classification problem since they all need no vacuum and overlap constraint (3.1) on the labeling functions in the unit simplex which could make their models to be non-convex. The second stage of our method is to find a binary partition via thresholding the smoothed result obtained from stage one. We proved that our proposed model has a unique solution and the derived primal-dual algorithm converges.

We tested our proposed method on four benchmark data sets and compared with the state-of-the-art methods. We also investigated the influence of the training sets selected uniformly and non-uniformly. For our method, different ways of generating initializations were implemented and validated. On the whole, the experimental results demonstrated that our method is superior in terms of classification accuracy and when parallel computing is considered, computation speed too. Therefore our method is an efficient and effective classification method for data sets like high-dimensional data or unstructured point clouds.

Acknowledgements. This work of R. Chan is partially supported by HKRGC Grants No. CityU12500915, CityU14306316, HKRGC CRF Grant C1007-15G, and HKRGC AoE Grant AoE/M-05/12. This work of T. Zeng is partially supported by the National Natural Science Foundation of China under Grant 11671002, CUHK start-up and CUHK DAG 4053296, 4053342. We thank Prof. Xue-Cheng Tai, Dr.

Ke Yin, Dr. Egil Bae and Prof. Ekaterina Merkurjev for providing the codes of their methods [1, 56].

REFERENCES

- [1] E. BAE AND E. MERKURJEV, *Convex variational methods on graphs for multiclass segmentation of high-dimensional data and point clouds*, Journal of Mathematical Imaging and Vision, 58 (2017), pp. 468–493.
- [2] L. BAR, T. F. CHAN, G. CHUNG, M. JUNG, N. KIRYATI, R. MOHIEDDINE, N. SOCHEN, AND L. A. VESE, *Mumford and Shah model and its applications to image segmentation and image restoration*, Springer New York, New York, NY, 2011, pp. 1095–1157.
- [3] B. BAUER, X. CAI, S. PETH, K. SCHLADITZ, AND G. STEIDL, *Variational-based segmentation of bio-pores in tomographic images*, Computers & Geosciences, 98 (2017), pp. 1–8.
- [4] A. BERTOZZI AND A. FLENNER, *Diffuse interface models on graphs for classification of high dimensional data*, Multiscale Modeling and Simulation, 10 (2012), pp. 1090–1118.
- [5] J. C. BEZDEK, R. EHRLICH, AND W. FULL, *FCM: The Fuzzy C-means clustering algorithm*, Computers & Geosciences, 10 (1984), pp. 191–203.
- [6] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and statistical learning via the alternating direction method of multipliers*, Foundations and Trends in Machine Learning, 3 (2011), pp. 1–122.
- [7] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, New York, NY, USA, 2004.
- [8] Y. BOYKOV AND G. FUNKA-LEA, *Graph cuts and efficient N-D image segmentation*, International Journal of Computer Vision, 70 (2006), pp. 109–131.
- [9] X. BRESSON, S. ESEDOGLU, P. VANDERGHEYNST, J. THIRAN, AND S. OSHER, *Fast global minimization of the active contour/snake model*, Journal of Mathematical Imaging and Vision, 28 (2007), pp. 151–167.
- [10] X. BRESSON, X.-C. TAI, T. F. CHAN, AND A. SZLAM, *Multi-class transductive learning based on l_1 relaxations of cheeger cut and mumford-shah-potts model*, Journal of Mathematical Imaging and Vision, 49 (2014), pp. 191–201.
- [11] E. BROWN, T. CHAN, AND X. BRESSON, *Completely convex formulation of the Chan-Vese image segmentation model*, International Journal of Computer Vision, 98 (2012), pp. 103–121.
- [12] N. BURNET, J. SCAIFE, M. ROMANCHIKOVA, S. THOMAS, AND ET AL., *Applying physical science techniques and CERN technology to an unsolved problem in radiation treatment for cancer: the multidisciplinary ‘VoxTox’ research programme*, CERN ideaSquare journal of experimental innovation, 1 (2017).
- [13] X. CAI, *Variational image segmentation model coupled with image restoration achievements*, Pattern Recognition, 48 (2015), pp. 2029–2042.
- [14] X. CAI, R. H. CHAN, S. MORIGI, AND F. SGALLARI, *Framelet-based algorithm for segmentation of tubular structures*, in Scale Space and Variational Methods in Computer Vision, A. M. Bruckstein, B. M. ter Haar Romeny, A. M. Bronstein, and M. M. Bronstein, eds., Berlin, Heidelberg, 2012, Springer Berlin Heidelberg, pp. 411–422.
- [15] X. CAI, R. H. CHAN, S. MORIGI, AND F. SGALLARI, *Vessel segmentation in medical imaging using a tight-frame-based algorithm*, SIAM Journal on Imaging Sciences, 6 (2013), pp. 464–486.
- [16] X. CAI, R. H. CHAN, M. NIKOLOVA, AND T. ZENG, *A three-stage approach for segmenting degraded color images: Smoothing, lifting and thresholding (SLaT)*, Journal of Scientific Computing, 72 (2017), pp. 1313–1332.
- [17] X. CAI, R. H. CHAN, C.-B. SCHÖNLIEB, G. STEIDL, AND T. ZENG, *Linkage between piecewise constant Mumford-Shah model and ROF model and its virtue in image segmentation*, arXiv:1807.10194, (2018).
- [18] X. CAI, R. H. CHAN, AND T. ZENG, *A two-stage image segmentation method using a convex variant of the Mumford-Shah model and thresholding*, SIAM Journal on Imaging Sciences, 6 (2013), pp. 368–390.
- [19] X. CAI, J. FITSCHEN, M. NIKOLOVA, G. STEIDL, AND M. STORATH, *Disparity and optical flow partitioning using extended Potts priors*, Information and Inference: A Journal of the IMA, 4 (2014), pp. 43–62.
- [20] X. CAI, C.-B. SCHÖNLIEB, J. LEE, AND ET AL., *Automatic contouring of soft organs for image-guided prostate radiotherapy*, Radiotherapy and Oncology, 119 (2016), pp. S895–S896.
- [21] X. CAI AND G. STEIDL, *Multiclass segmentation by iterated ROF thresholding*, in Energy Minimization Methods in Computer Vision and Pattern Recognition, A. Heyden, F. Kahl,

- C. Olsson, M. Oskarsson, and X.-C. Tai, eds., Berlin, Heidelberg, 2013, Springer Berlin Heidelberg, pp. 237–250.
- [22] A. CHAMBOLLE, M. NOVAGA, D. CREMERS, AND T. POCK, *An introduction to total variation for image analysis*, in in Theoretical Foundations and Numerical Methods for Sparse Recovery, De Gruyter, 2010.
- [23] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications to imaging*, Journal of Mathematical Imaging and Vision, 40 (2011), pp. 120–145.
- [24] R. H. CHAN AND M. NG, *Conjugate gradient method for Toeplitz systems*, SIAM Review, 38 (1996), pp. 427–482.
- [25] R. H. CHAN, H. YANG, AND T. ZENG, *A two-stage image segmentation method for blurry images with Poisson or multiplicative Gamma noise*, SIAM Journal on Imaging Sciences, 7 (2014), pp. 98–127.
- [26] T. F. CHAN, S. ESEDOGLU, AND M. NIKOLOVA, *Algorithms for finding global minimizers of image segmentation and denoising models*, SIAM Journal on Applied Mathematics, 66 (2006), pp. 1632–1648.
- [27] T. F. CHAN AND L. A. VESE, *Active contours without edges*, IEEE Transactions on Image Processing, 10 (2001), pp. 266–277.
- [28] O. CHAPELLE, B. SCHLKOPF, AND A. ZIEN, *Semi-Supervised Learning*, The MIT Press, 1st ed., 2010.
- [29] C. CORTES AND V. VAPNIK, *Support-vector networks*, Machine Learning, 20 (1995), pp. 273–297.
- [30] D. CREMERS, M. ROUSSON, AND R. DERICHE, *A review of statistical approaches to level set segmentation: integrating color, texture, motion and shape*, International Journal of Computer Vision, 72 (2007), pp. 195–215.
- [31] B. DONG, A. CHIEN, AND Z. SHEN, *Frame based segmentation for medical images*, Communications in Mathematical Sciences, 9 (2010), pp. 551–559.
- [32] A. ELMOATAZ, O. LEZORAY, AND S. BOUGLEUX, *Nonlocal discrete regularization on weighted graphs: a framework for image and manifold processing*, IEEE Transactions on Image Processing, 17 (2008), pp. 1047–1060.
- [33] C. GARCIA-CARDONA, E. MERKURJEV, A. L. BERTOZZI, A. FLENNER, AND A. G. PERCUS, *Multiclass data segmentation using diffuse interface methods on graphs*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 36 (2014), pp. 1600–1613.
- [34] T. GOLDSTEIN AND S. OSHER, *The split Bregman method for l_1 -regularized problems*, SIAM Journal on Imaging Sciences, 2 (2009), pp. 323–343.
- [35] Y. HE, M. Y. HUSSAINI, J. MA, B. SHAFEI, AND G. STEIDL, *A new Fuzzy C-means method with total variation regularization for image segmentation of images with noisy and incomplete data*, Pattern Recognition, 45 (2012), pp. 3463–3471.
- [36] M. HEIN AND S. SETZER, *Beyond spectral clustering - tight relaxations of balanced graph cuts*, in Advances in Neural Information Processing Systems 24, J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, eds., Curran Associates, Inc., 2011, pp. 2366–2374.
- [37] T. KANUNGO, D. M. MOUNT, N. S. NETANYAHU, C. D. PIATKO, R. SILVERMAN, AND A. Y. WU, *An efficient k-means clustering algorithm: analysis and implementation*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 24 (2002), pp. 881–892.
- [38] J. LEE, X. CAI, J. LELLMANN, M. DALPONTE, AND ET AL., *Individual tree species classification from airborne multisensor imagery using robust PCA*, IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing, 9 (2016), pp. 2554–2567.
- [39] J. LELLMANN AND C. SCHNÖRR, *Continuous multiclass labeling approaches and algorithms*, SIAM Journal on Imaging Sciences, 44 (2011), pp. 1049–1096.
- [40] O. LÉZORAY, A. ELMOATAZ, AND V. T. TA, *Nonlocal PDEs on graphs for active contours models with applications to image segmentation and data clustering*, IEEE International Conference on Acoustics, Speech and Signal Processing, (2012), pp. 873–876.
- [41] F. LI, M. NG, T. ZENG, AND C. SHEN, *A multiphase image segmentation method based on fuzzy region competition*, SIAM Journal on Imaging Sciences, 3 (2010), pp. 277–299.
- [42] E. MERKURJEV, E. BAE, A. L. BERTOZZI, AND X.-C. TAI, *Global binary optimization on graphs for classification of high-dimensional data*, Journal of Mathematical Imaging and Vision, 52 (2015), pp. 414–435.
- [43] E. MERKURJEV, A. BERTOZZI, X. YAN, AND K. LERMAN, *Modified Cheeger and ratio cut methods using the Ginzburg-Landau functional for classification of high-dimensional data*, Inverse Problems, 33 (2017), p. 074003.
- [44] E. MERKURJEV, T. KOSTIC, AND A. BERTOZZI, *An MBO scheme on graphs for classification*

- and image processing, SIAM Journal on Imaging Sciences, 6 (2013), pp. 1903–1930.
- [45] B. MERRIMAN AND S. J. RUUTH, *Diffusion generated motion of curves on surfaces*, Journal of Computational Physics, 225 (2007), pp. 2267–2282.
 - [46] D. MUMFORD AND J. SHAH, *Optimal approximations by piecewise smooth functions and associated variational problems*, Communications on Pure and Applied Mathematics, 42 (1989), pp. 577–685.
 - [47] B. OSTING, C. WHITE, AND E. OUDET, *Minimal Dirichlet energy partitions for graphs*, SIAM Journal on Imaging Sciences, 36 (2014), pp. 1635–1651.
 - [48] T. POCK, A. CHAMBOLLE, D. CREMERS, AND H. BISCHOF, *A convex relaxation approach for computing minimal partitions*, IEEE Conference on Computer Vision and Pattern Recognition, (2009), pp. 810–817.
 - [49] T. POCK, D. CREMERS, H. BISCHOF, AND A. CHAMBOLLE, *An algorithm for minimizing the mumford-shah functional*, in 2009 IEEE 12th International Conference on Computer Vision, IEEE, 2009, pp. 1133–1140.
 - [50] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Physica D: Nonlinear Phenomena, 60 (1992), pp. 259–268.
 - [51] J. SCAIFE, K. HARRISON, A. DREW, AND ET AL., *Accuracy of manual and automated rectal contours using helical tomotherapy image guidance scans during prostate radiotherapy*, Journal of Clinical Oncology, 33 (2015), p. 94.
 - [52] J. SHI AND J. MALIK, *Normalized cuts and image segmentation*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 22 (2000), pp. 888–905.
 - [53] C. SILPA-ANAN AND R. HARTLEY, *Optimised KD-trees for fast image descriptor matching*, IEEE Conference on Computer Vision and Pattern Recognition, (2008), pp. 1–8.
 - [54] A. SUBRAMANYA AND J. BILMES, *Semi-supervised learning with measure propagation*, Journal of Machine Learning Research, 12 (2011), pp. 3311–3370.
 - [55] L. VESE AND T. F. CHAN, *A multiphase level set framework for image segmentation using the Mumford and Shah model*, International Journal of Computer Vision, 50 (2002), pp. 271–293.
 - [56] K. YIN AND X.-C. TAI, *An effective region force for some variational models for learning and clustering*, Journal of Scientific Computing, 74 (2018), pp. 175–196.
 - [57] S. YU AND J. SHI, *Multiclass spectral clustering*, in Proceedings Ninth IEEE International Conference on Computer Vision, Oct 2003, pp. 313–319 vol.1.
 - [58] J. YUAN, E. BAE, X.-C. TAI, AND Y. BOYKOV, *A continuous max-flow approach to potts model*, in European Conference on Computer Vision, 2010, pp. 379–392.
 - [59] Y. ZHANG, B. MATUSZEWSKI, L. SHARK, AND C. MOORE, *Medical image segmentation using new hybrid level-set method*, in 2008 Fifth International Conference BioMedical Visualization: Information Visualization in Medical and Biomedical Informatics, 2008, pp. 71–76.