

Uncertainty Quantification for Semi-supervised Multi-class Classification in Image Processing and Ego-Motion Analysis of Body-Worn Videos

Yiling Qiao¹, Chang Shi², Chenjian Wang³, Hao Li³, Matt Haberland⁴, Xiyang Luo³, Andrew M. Stuart⁵, Andrea L. Bertozzi³;

¹ University of Chinese Academy of Science; Beijing, China

² Renmin University of China; Beijing, China

³ University of California, Los Angeles; Los Angeles, California, USA

⁴ California Polytechnic State University; San Luis Obispo, California, USA

⁵ California Institute of Technology; Pasadena, California, USA

Abstract

Semi-supervised learning uses underlying relationships in data with a scarcity of ground-truth labels. In this paper, we introduce an uncertainty quantification (UQ) method for graph-based semi-supervised multi-class classification problems. We not only predict the class label for each data point, but also provide a confidence score for the prediction. We adopt a Bayesian approach and propose a graphical multi-class probit model together with an effective Gibbs sampling procedure. Furthermore, we propose a confidence measure for each data point that correlates with the classification performance. We use the empirical properties of the proposed confidence measure to guide the design of a human-in-the-loop system. The uncertainty quantification algorithm and the human-in-the-loop system are successfully applied to classification problems in image processing and ego-motion analysis of body-worn videos.

Introduction

Applications such as police body-worn video cameras generate a huge amount of data, beyond what is humanly possible for analysts to review. Such problems are ripe for the development of semi-supervised learning algorithms, which, by definition, use a small amount of training data. In the last year, progress has been made in applying graph-based semi-supervised learning to body-worn videos with the goal of recognizing camera-wearers' activities, i.e., ego-motion [11, 5]. However, as is often the case with real-world videos, the variability of the data leads to imperfect classification. Recently, the authors of [3] proposed to pair uncertainty quantification (UQ) with the binary classification problem on a similarity graph. Besides a label assigned to each data point, they also estimated a measure of uncertainty, which helped identify hard-to-classify data points that require further investigation.

In the present paper, we push the UQ methodology to a multi-class setting. We extend the binary graphical probit method to a multi-class version and develop a Gibbs sampler that draws samples from the posterior distribution. We propose a confidence measure for each data point that we find correlates with the classification performance; we observe that data points with higher confidence scores are more likely to be classified correctly. Along with the new methodology and the empirical observations, we develop the foundations for a system with a human in the loop who serves to provide additional class labels based on the confidence

scores; our uncertainty quantification method identifies hard-to-classify data points and the human in the loop assigns ground truth to them, leading to reduced overall confidence scores. Our ideas are tested on an image data set — the MNIST data set [10] — and a body-worn video data set, the HUJI EgoSeg data set [14].

Related Work

Semi-supervised learning has been studied extensively in the past two decades and has been successfully applied to applications such as hyperspectral images [12] and body-worn videos [11, 5]. We refer readers to [22] and the more recent article [1] for a literature review. We focus on graph-based methods, in which a similarity is measured for each pair of nodes (i.e. data points) and label information is spread across the similarity graph from a small set of labeled fidelity points. The similarity information is often leveraged via the graph Laplacian, which has been used in a myriad of machine learning methods (see, for instance, [19, 20, 21, 23]). The analogy between the graph Laplacian and the classical Laplacian operator inspires a number of PDE-based classification methods, such as [2, 9]; this also introduces the recent development in uncertainty quantification to the machine learning community. For instance, in their recent work [3], the authors used an efficient sampling method that was originally developed for PDE-based inverse problems [6] to perform uncertainty quantification for binary classification.

We refer readers to the books [16, 17] and the recent article [13] for a review of methodologies employed in the field of uncertainty quantification. For the specific application to machine learning methods, the book [20] investigates uncertainty quantification for a variety machine learning problems using a Gaussian process prior. Except the above-mentioned book and the recent work [3], most machine learning methods, even those developed with the Bayesian way of thinking, focus on finding the optimal classification (and/or hyperparameters that produce the optimal classification) in an optimization context and do not consider or utilize uncertainty quantification.

Methodology

Graphical Setting

Let $X = \{x_1, x_2, \dots, x_n\}$ be a set of feature vectors, where $x_i \in \mathbb{R}^d$. Let $Z = \{1, 2, \dots, n\}$ index the entire dataset and $Z' \subset Z$

be a fidelity set consisting of nodes with known labels. We aim to classify n data points into c classes such that:

- 1) data points with similar feature vectors, measured via a suitable similarity measure, should belong to the same class;
- 2) the classification should respect the ground-truth labels on the fidelity set.

We consider each data point as a node in a weighted similarity graph, where the edge weights are given by

$$w_{ij} = \exp\left(-\|x_i - x_j\|^2 / \tau_{ij}\right),$$

where $\|\cdot\|$ is the Euclidean distance and τ_{ij} are the self-tuning constants proposed in [21]. The weights are chosen such that a pair of nodes with similar feature vectors will have a weight close to one and dissimilar nodes will have a near-zero weight. Suppose $u: Z \rightarrow \mathbb{R}^c$ is an assignment function; if $u_\ell(i) = \max_{\ell} u_\ell(i)$ then we interpret this to mean that u assigns class ℓ to data point i . One way to achieve a classification is to optimize the following objective function with respect to an assignment function u :

$$J(u) = \frac{1}{4} \sum_{i,j=1}^n w_{ij} \|u(i) - u(j)\|^2 + \Phi(u, u'), \quad (1)$$

where u' encodes the ground-truth labels on the fidelity set Z' , $\Phi(u, u')$ measures the extent to which u differs from u' on Z' . Minimizing the first term in the object function ensures that a pair of data points (i, j) with a high similarity w_{ij} will be assigned to the same class.

Using matrix notation, we identify u and u' with $n \times c$ matrices so that $u_{i\ell} = u_\ell(i)$. If we let W be the matrix of w_{ij} and $D = \text{diag}(d_1, d_2, \dots, d_n)$ where $d_i = \sum_j w_{ij}$, we can introduce the graph Laplacian

$$L = D - W \quad (2)$$

and the Dirichlet energy

$$\langle u, Lu \rangle = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \|u(i) - u(j)\|^2, \quad (3)$$

where $\langle u, v \rangle = \text{trace}(u^T v)$, and hence we may write eq. (1) as

$$J(u) = \frac{1}{2} \langle u, Lu \rangle + \Phi(u, u') \quad (4)$$

The quadratic form in eq. (3) alludes to the connection to Bayesian Gaussian process models.

It is common in graph-based learning methods to use normalized variants of the graph Laplacian in place of the unnormalized graph Laplacian eq. (2) because of better numerical properties as well as the classification performance (see, for instance, [2]). One popular choice is the symmetrically normalized graph Laplacian,

$$L_{\text{sym}} = D^{-1/2} L D^{-1/2}, \quad (5)$$

which is convenient to compute with due to its symmetry. With the choice of the symmetrically normalized graph Laplacian, the quadratic form in eq. (3) becomes

$$\langle u, L_{\text{sym}} u \rangle = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \left\| \frac{u(i)}{\sqrt{d_i}} - \frac{u(j)}{\sqrt{d_j}} \right\|^2.$$

In the remainder of this manuscript, the notation L is a placeholder for any choice of graph Laplacian.

Bayesian model

We now present a Bayesian model for the assignment function u , for which the posterior distribution takes the form:

$$p(u|u') \propto \exp(-J(u)), \quad (6)$$

so a maximum a posteriori probability (MAP) estimator is a minimizer of $J(u)$. We assume the prior on u is a Gaussian distribution,

$$p(u) \propto \exp\left(-\frac{1}{2} \langle u, Lu \rangle\right).$$

To explicitly construct a sample u that follows the prior distribution, we employ the Karhunen-Loève expansion. Let $L = Q\Lambda Q^T$ be the eigen-decomposition of the graph Laplacian where the columns of $Q \in \mathbb{R}^{n \times n}$ form an orthonormal basis of \mathbb{R}^n and $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ obeys

$$0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n.$$

We observe that L is positive semi-definite. Suppose $\{\xi_i\}_{i=1}^n$ is a collection of independent c -variate normal random variables $\mathcal{N}(0, I_c)$, where I_c is an identity matrix of size c . We construct a sample u as the random sum

$$u = \sum_{i=2}^n \lambda_i^{-1/2} q_i \xi_i^T,$$

so that the columns of u live in $\text{span}\{q_1\}^\perp$ and u has the desired probability distribution

$$p(u) \propto \exp\left(-\frac{1}{2} \sum_{i=1}^n \sum_{\ell=1}^c \lambda_i \langle u_\ell, q_i \rangle^2\right) = \exp\left(-\frac{1}{2} \langle u, Lu \rangle\right). \quad (7)$$

In graph-based semi-supervised learning, it is common to approximate the graph Laplacian using only its first few eigenvectors and eigenvalues, since these often contain the relevant geometric information reflecting clustering of the data points at nodes of the graph (see, for instance, [19]). Such truncation of the spectrum both reduces computation cost and often improves the classification performance. In this case, we also employ spectral truncation and let u be the random sum up to K for $K \ll n$, i.e.

$$u = \sum_{i=2}^K \lambda_i^{-1/2} q_i \xi_i^T. \quad (8)$$

In their recent work [3], the authors considered several likelihood functions $p(u'|u)$ to connect the latent variable u to the ground-truth labeling u' for binary classification. In the present paper, we primarily investigate the independent probit likelihood function. Suppose $\{\eta(i)\}_{i \in Z'}$ is a collection of independent c -variate normal random variables $\mathcal{N}(0, \gamma^2 I_c)$ where γ^2 is the noise variance. We connect u to u' via

$$\begin{aligned} v(i) &= u(i) + \eta(i) \\ u'(i) &= \text{threshold}(v(i)), i \in Z'. \end{aligned}$$

The threshold operator applied to a vector simply sets the largest element in the vector to be 1 and the rest to be 0. With the introduction of latent variables $\{v(i)\}_{i \in Z'}$, we have, from Bayes formula, the following joint posterior probability distribution

$$p(u, v|u') \propto \exp\left(-\frac{1}{2}\langle u, Lu \rangle - \frac{1}{2\gamma^2} \sum_{i \in Z'} \|u(i) - v(i)\|^2\right) \times \prod_{i \in Z'} 1_{\text{threshold}(v(i))=u'(i)}.$$

Using the change of variable from u to ξ , for $\xi = (\xi_1, \xi_2, \dots, \xi_K) \in \mathbb{R}^{c \times K}$ in eq. (8), we can apply our chosen sampling method to $p(\xi, v|u')$. We compute that the joint probability

$$p(\xi, v|u') \propto \exp\left(-\frac{1}{2}\langle \xi^T, \Lambda' \xi^T \rangle - \frac{1}{2\gamma^2} \|HQ' \xi^T - v\|^2\right) \times \prod_{i \in Z'} 1_{\text{threshold}(v(i))=u'(i)},$$

where $\Lambda' = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_K)$, the matrix $Q' \in \mathbb{R}^{n \times K}$ consists of the first K eigenvectors of the graph Laplacian, and $H = (\delta_{ij}) \in \mathbb{R}^{|Z'| \times n}$ for $Z' = \{j_i : i = 1, 2, \dots, |Z'|\}$. We note that H applied to a matrix selects its rows of the fidelity set Z' .

To sample from the joint posterior distribution, a Gibbs sampler will alternate between the following three steps:

- 1) Draw ξ from $p(\xi|v, u')$,
- 2) Construct u from ξ via eq. (8),
- 3) Draw v from $p(v|u, u')$.

For Step 1) we note that for each $\ell \in \{1, 2, \dots, c\}$, the conditional probability for each row of ξ , denoted as $p(\xi_{:, \ell}|v, u')$ has the same distribution as

$$\mathcal{N}(m, P^{-1}), P = \Lambda' + \frac{1}{\gamma^2} Q'^T H^T H Q', m = \frac{1}{\gamma^2} P^{-1} Q'^T H^T v_\ell.$$

In Step 3), for each $i \in Z'$, we need to sample a c -variate normal random variable subject to a linear inequality constraint; let a_i denote the unique index such that $u'_{a_i}(i) = 1$ for $i \in Z'$, i.e., data point i belongs to class a_i according to the ground-truth label. Then we need to sample $v(i)$ according to the following conditions:

$$v(i) \sim \mathcal{N}(u(i), \gamma^2 I_c), \quad v_{a_i}(i) \geq v_\ell(i) \text{ for all } \ell \in \{1, 2, \dots, c\}.$$

We use the implementation from [4] to efficiently draw samples from the linearly constrained normal distribution.

Uncertainty Quantification

Given a set of samples $\{u^{(k)}\}_{k=1}^N$ from the Gibbs sampler, we investigate $\mathbb{E}_{u|u'}(\text{threshold}(u))$, the posterior mean of $\text{threshold}(u)$; this can be approximated by the sample mean

$$s_\ell(i) = \mathbb{E}_{u|u'}(\text{threshold}(u(i)))_\ell \approx \frac{1}{N} \sum_{k=1}^N \text{threshold}(u(i))_\ell.$$

Since each element $\text{threshold}(u(i))_\ell$ is either zero or one, the expectation $s_\ell(i)$ simply gives the probability, under the posterior distribution, of the element being one; that is $s_\ell(i)$ can be interpreted as the probability data point i belongs to class ℓ . We note

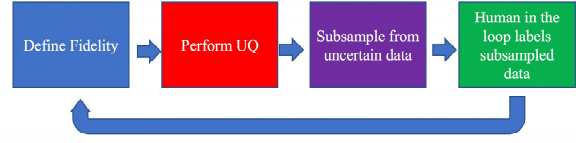


Figure 1. A flow chart summarizing the proposed human-in-the-loop system.

that for each data point, the probability of it belonging to each class should sum to one, i.e., $\sum_\ell s_\ell(i) = 1$. This is obeyed by both the posterior mean and the sample mean approximation. We can use the posterior mean $s(i)$ as a classifier, which classifies data point i according to its largest entry.

Intuitively, a single large $s_\ell(i)$ for a data point i indicates a very confident classification of class ℓ ; in this case, the remaining entries in $s(i)$ are necessarily small due to the sum-to-one condition; this creates a large variance in the vector $s(i)$. On the other hand, if entries in the vector $s(i)$ are all roughly equal, meaning the data point is equally likely to be classified as either class, the classification has a lot of uncertainty, resulting in an $s(i)$ with a small variance. Based on this intuition, we measure the classification confidence of node i by the variance of $s(i)$

$$S(i) = \text{var}(s(i)) = \frac{1}{c} \sum_{\ell=1}^c \left(s_\ell(i) - \frac{1}{c} \sum_{\ell=1}^c s_\ell(i) \right)^2.$$

We emphasize that this variance is not the posterior variance. However, we can show the following connection between the quantity $S(i)$ and the posterior variance

$$S(i) = \frac{1}{c} - \frac{1}{c^2} - \frac{1}{c} \sum_{\ell=1}^c \text{var}_{u|u'}(\text{threshold}(u(i)))_\ell,$$

where $\text{var}_{u|u'}(\cdot)$ is the posterior variance. Therefore, the quantity $S(i)$ is a constant minus the mean posterior variance, which can be interpreted as a measure of uncertainty, averaged over all classes.

Human-in-the-loop

In the following experiments section, we demonstrate a positive correlation between the proposed confidence score and the classification performance; the confidence score enables us to locate hard-to-classify data points, which we may label and use as additional fidelity. This naturally leads to the idea of using the confidence measure to intelligently select new fidelity points to achieve a better classification performance with limited human labeling effort. We design a human-in-the-loop system as follows (see fig. 1). We start with a small set of initial fidelity points and apply the UQ algorithm to obtain a confidence score for the entire data set. We randomly sample, in a uniform fashion, additional candidate fidelity points with confidence scores within a percentile range. The human in the loop then observes each of the candidate fidelity points to assign ground truth to them. We perform the UQ algorithm again to update the confidence scores and repeat the process until we reach the maximum number of fidelity points permitted (this will be determined by the application).

We observe that in practice adding data points with the lowest confidence scores does not benefit overall classification per-

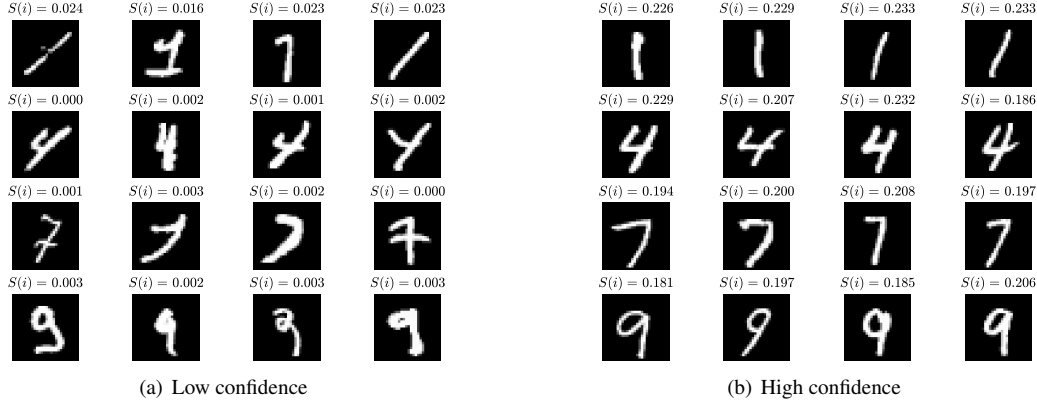


Figure 2. Uncertainty quantification on the MNIST dataset. $S(i)$ is the proposed confidence score. For each digit, we present four examples chosen from the top/bottom ten with the highest/lowest confidence scores within each class.

formance because these data points are often outliers. The significance of classifying these outliers correctly is scenario dependent. In our experiments, we focus on the overall accuracy and do not sample fidelity from data points with confidence scores strictly below the tenth percentile.

Experiments

We perform uncertainty quantification on 1) the MNIST data set, a handwritten digit data set, and 2) the HUJI EgoSeg data set, a body-worn video data set. Through these experiments, we illustrate some empirical properties of the confidence score; we demonstrate its correlation with the classification performance. We also validate our human-in-the-loop system and showcase its ability to improve classification results with limited human input.

MNIST

The MNIST data set [10] consists of 70,000 images of handwritten digits; each image is of the size 28×28 pixels. We choose uniformly at random 500 images each from the digits 1, 4, 7, and 9 to form a graph of 2000 nodes. We follow the graph construction procedure in [3]; each image is projected onto the lead 50 principal components yielding a 50-dimensional feature vector, and we construct a 15-nearest neighbor graph. The weighting constants τ_{ij} are chosen according to [21]. For data point i , we compute the mean distance of its 15 nearest neighbors, denoted as τ_i ; then the weighting constant τ_{ij} is given by $\tau_{ij} = \tau_i \tau_j$. We use the symmetrically normalized graph Laplacian (see eq. (5)) and truncate its spectrum at $K = 300$. We perform the Gibbs sampler with 3% uniformly randomly sampled fidelity points; the noise variance is chosen to be $\gamma = 0.1$, and we draw 2×10^4 samples to estimate the uncertainty. We showcase examples of images with the highest or the lowest confidence scores in fig. 2. It is interesting to note that the lowest confidence score of digit 1 is much higher than that of the other digits; we theorize that it is easier for the algorithm to differentiate digit 1 from the other three digits.

Body-Worn Videos

We also apply our method to the HUJI EgoSeg data set [14, 15]. This data set contains 65 hours of egocentric videos including 44 videos filmed using a head-mounted GoPro Hero3+,

the Disney data set [7] and other YouTube videos¹. In the recent paper [5], a graph-based semi-supervised learning method is applied to this data set to classify video segments according to camera-wearers' activities and showed promising results. This data set consists of footage of 7 activities: *Walking, Driving, Riding Bus, Biking, Standing, Sitting, and Static*. We follow the same feature extraction procedure described in [5] to obtain a 50-dimensional feature vector for every 4-second video segments; this yields 36,421 segments. To speed up our calculation, we sample every fifth segment. The graph is constructed from the 50-dimensional feature vectors, and the weighting constants $\tau_{ij} = \tau_i \tau_j$ are chosen according to [21], where τ_i is the distance of the 40th nearest neighbor of node i . We employ the symmetrically normalized graph Laplacian and truncate the spectrum at $K = 400$. The eigenvectors are computed using a low-rank approximation of the graph Laplacian via the Nystrom extension [8]. The Gibbs sampler is applied with $\gamma = 0.1$ and 2×10^4 iterations.

The data set is separated into a training and testing set, which are disjoint sets of videos; the training set contains around 65% of the data, measured in terms of the footage length. We refer readers to [15] for the details of the experimental protocol. However, we do not use the full training set but instead take a portion of it as the fidelity. All classification performances are evaluated on the testing set only. We first investigate the correlation between the confidence score and the classification accuracy. We perform uncertainty quantification with 12% of the training set. Recall that the classification is produced by taking the largest entry of the posterior mean $s(i)$ for each data point i . In fig. 3, we plot the classification accuracy of the top $x\%$ to $x+5\%$ confident data points for each $x \in \{0, 5, 10, \dots, 95\}$. We observe that the classification is more accurate on data points with higher confidence scores. We also validate our human-in-the-loop system on this data set. We start with 6% fidelity data and gradually increase the fidelity percentage to 30% over five iterations; at each iteration, we introduce additional 6% fidelity points sampled from data points with confidence scores within in the range of the tenth and 50th percentile. We perform uncertainty quantification as well as a graph-based semi-supervised learning method (an MBO scheme [2]) using the same set of fidelity points. We refer readers to the appendix for a description of the MBO scheme and its

¹<http://www.vision.huji.ac.il/egoseg/>

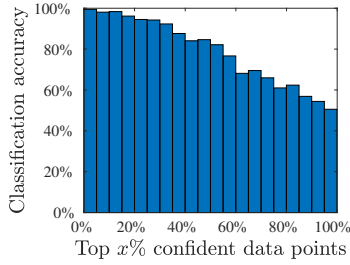
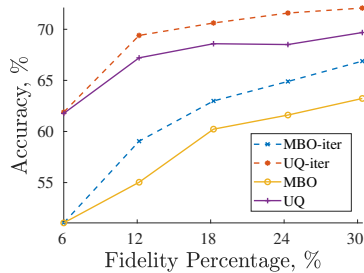
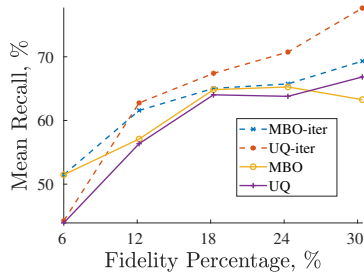


Figure 3. Classification accuracy on data points with top $x\%$ to $x+5\%$ confidence scores on the HUJI EgoSeg data set. We group data points based on their confidence score; each group contains 5% data points and we evaluate the classification accuracy on each group.



(a) Accuracy



(b) Recall

Figure 4. Classification performance of UQ and an MBO classifier using iteratively generated fidelity (UQ/MBO-iter) and uniformly randomly sampled fidelity (UQ/MBO) on the HUJI EgoSeg data set.

parameters that we use for this experiment. We compare the classification performance, measured in terms of accuracy and mean recall averaged over seven classes, of both classifiers using iteratively generated fidelity against the same classifiers using uniformly randomly sampled fidelity. The results are presented in fig. 4. We observe that both classifiers benefit from the intelligently sampled fidelity in terms of producing higher accuracy and mean recall than using uniformly randomly sampled fidelity.

Conclusion

In this paper, we considered the problem of uncertainty quantification in a graph-based semi-supervised multi-class classification problem. We extended the graphical probit model, originally proposed for the binary classification problem in [3], to the multi-class case. We proposed a Gibbs sampler to sample from the posterior distribution and a confidence score that connects to the posterior variance. Through our experiments on the MNIST

data set, we demonstrated that the proposed confidence score is easy to interpret; it is clear to see the contrast between the digit images with low confidence scores and ones with high confidence scores. The proposed confidence score also exhibits a correlation with the classification performance in our experiments on the HUJI EgoSeg data set. Based on these observations, we designed a human-in-the-loop system to efficiently use human labeling effort to improve classification results. We validated this system on the HUJI EgoSeg data set and observed that the classifiers that we studied produced improved classification using the human-in-the-loop system than the same classifiers using uniformly randomly sampled fidelity.

Moving forward, we can develop new theory of uncertainty quantification for semi-supervised multi-class classification. We can investigate the performance bound of the Gibbs sampler with respect to a large number of data points and classes. We can extend the previous analysis of uncertainty quantification methods for binary classification to the multi-class case, in which we suspect the number of classes play a nontrivial role in the performance of the sampling methods. We also point out that speed is the primary concern of the current Gibbs sampler. Despite the development of scalable graph-based semi-supervised learning methods (see [2] for an example), the Gibbs sampler is mostly sequential; we draw each sample based on the previous one. Nevertheless, the current work opens the door for the development of a system that combines modern machine learning with expert analyst knowledge.

Acknowledgement

We thank Hannah Droege, Sara Tro, and Yang Wang for useful comments. We acknowledge support from NSF grants DMS-1737770 and DMS-1417674. Yiling Qiao and Chang Shi were supported by the UCLA-CSST program. AMS is funded by NSF grant DMS 1818977.

References

- [1] A. L. Bertozzi. Graphical models in machine learning, networks and uncertainty quantification. In *proceedings of 2018 International Congress of Mathematicians*, volume 3, pages 3853–3880.
- [2] A. L. Bertozzi and A. Flenner. Diffuse interface models on graphs for classification of high dimensional data. *SIAM Review*, 58(2):293–328, 2016.
- [3] A. L. Bertozzi, X. Luo, A. M. Stuart, and K. C. Zygalakis. Uncertainty quantification in the classification of high dimensional data. *SIAM/ASA J. Uncertainty Quantification*, 6(2):568–595, 2018.
- [4] Z. I. Botev and P. L’Ecuyer. Efficient probability estimation and simulation of the truncated multivariate student-t distribution. In *Proceedings of the 2015 Winter Simulation Conference*, pages 380–391. IEEE Press, 2015.
- [5] H. Chen, H. Li, A. Song, M. Haberland, O. Akar, A. Dhillon, T. Zhou, A. L. Bertozzi, and J. P. Brantingham. First-person activity recognition in body-worn video. 2018.
- [6] S. L. Cotter, G. O. Roberts, A. M. Stuart, and D. White. MCMC methods for functions: modifying old algorithms to make them faster. *Statistical Science*, pages 424–446, 2013.
- [7] A. Fathi, J. K. Hodgins, and J. M. Rehg. Social interactions: A first-person perspective. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 1226–1233. IEEE, 2012.
- [8] C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral grouping

using the nystrom method. *IEEE transactions on pattern analysis and machine intelligence*, 26(2):214–225, 2004.

- [9] 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(8):1600–1613, 2014.
- [10] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [11] Z. Meng, J. Sanchez, J.-M. Morel, A. Bertozzi, and B. P. Jeffrey. Ego-motion classification for body-worn videos. In *2016 Conference on Imaging, vision and learning based on optimization and PDEs*, 2016.
- [12] E. Merkurjev, J. Sunu, and A. L. Bertozzi. Graph MBO method for multiclass segmentation of hyperspectral stand-off detection video. In *2014 IEEE International Conference on Image Processing (ICIP)*.
- [13] H. Owghadi, C. Scovel, T. J. Sullivan, M. McKerns, and M. Ortiz. Optimal uncertainty quantification. *Siam Review*, 55(2):271–345, 2013.
- [14] Y. Poleg, C. Arora, and S. Peleg. Temporal segmentation of egocentric videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2537–2544, 2014.
- [15] Y. Poleg, A. Ephrat, S. Peleg, and C. Arora. Compact CNN for indexing egocentric videos. In *Applications of Computer Vision (WACV), 2016 IEEE Winter Conference on*, pages 1–9. IEEE, 2016.
- [16] R. C. Smith. *Uncertainty quantification: theory, implementation, and applications*, volume 12. SIAM, 2013.
- [17] T. J. Sullivan. *Introduction to uncertainty quantification*, volume 63. Springer, 2015.
- [18] Y. Van Gennip, N. Guillen, B. Osting, and A. L. Bertozzi. Mean curvature, threshold dynamics, and phase field theory on finite graphs. *Milan Journal of Mathematics*, 82(1):3–65, 2014.
- [19] U. Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.
- [20] C. K. Williams and C. E. Rasmussen. Gaussian processes for regression. In *Advances in neural information processing systems*, pages 514–520, 1996.
- [21] L. Zelnik-Manor and P. Perona. Self-tuning spectral clustering. In *Advances in neural information processing systems*, pages 1601–1608, 2005.
- [22] X. Zhu. Semi-supervised learning literature survey. *Computer Science, University of Wisconsin-Madison*, 2(3):4, 2006.
- [23] X. Zhu, Z. Ghahramani, and J. D. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003.

Appendix

We detail the Merriman-Bence-Osher (MBO) scheme, a graph-based semi-supervised learning method. We optimize the graph Total Variation (TV) plus a least-squares fidelity term

$$\frac{1}{2}|u|_{\text{TV}} + \Phi_{\text{LS}}(u, u'). \quad (9)$$

subject to the constraint that each $u(i)$ is discrete; it lies on the corners of a unit simplex, i.e., one and only one entry of each $u(i)$ is one and the rest are zero; the graph TV is given by

$$|u|_{\text{TV}} = \frac{1}{2} \sum_{i,j=1}^n w_{ij} \|u(i) - u(j)\|_1, \quad (10)$$

and the least-squares fidelity term takes the form:

$$\Phi_{\text{LS}}(u, u') = \sum_{i \in \mathcal{Z}'} \frac{1}{2\gamma^2} \|u(i) - u'(i)\|^2.$$

We note that when $u(i)$ is discrete, the graph TV in eq. (10) agrees with the Dirichlet energy eq. (3). We then relax the combinatorial optimization problem; we allow $u(i)$ to take values in \mathbb{R}^c and penalize u for being away from the corners of the unit simplex with a multi-well potential

$$M(u) = \sum_{i=1}^n \prod_{\ell=1}^c \frac{1}{4} \|u(i) - e_\ell\|^2,$$

where e_ℓ is the unit vector in \mathbb{R}^c in the ℓ th direction. Replacing the graph TV with the Dirichlet energy and the discrete constraint with the addition of the multi-well potential, we arrive at the following objective function

$$\frac{1}{2} \langle u, Lu \rangle + \frac{1}{\varepsilon} M(u) + \Phi_{\text{LS}}(u, u') \quad (11)$$

for a small positive constant ε . The first two terms of eq. (11) is known as the Ginzburg-Landau functional, which Γ -converges to the graph TV as $\varepsilon \rightarrow 0$ [18].

In the MBO scheme, we alternatively perform the following two steps to update u :

1. *Diffuse*. Solve a force-driven heat equation

$$\frac{\partial u}{\partial t} = -Lu - \frac{1}{\gamma^2} \sum_{i \in \mathcal{Z}'} u(i) - u'(i),$$

for a short time Δt to obtain u^* ; this is effectively a gradient descent step for the first and third term of the objective function eq. (11).

2. *Threshold*. set

$$u(i) = e_\ell(i), \ell = \arg \max_{\ell} u_\ell^*(i).$$

This approximates the gradient descent step for the second term of eq. (11) when ε is small.

We use a semi-implicit method to solve the heat equation:

$$\frac{u^+ - u}{\delta t} = -Lu^+ - \frac{1}{\gamma^2} \sum_{i \in \mathcal{Z}'} u(i) - u'(i),$$

where $\delta t = \Delta t / N_{\text{step}}$ and N_{step} is the number of time steps used to solve the heat equation. We note that we use an implicit stepping for the term involving the graph Laplacian to resolve the potential stiffness of L . To accelerate the computation of the implicit stepping, we truncate the spectral at some level K , i.e., we approximate L by

$$Q' \Lambda' Q'^T = \sum_{i=1}^K \lambda_i q_i q_i^T.$$

In the experiment on the HUJI EgoSeg data set, we use the following set of parameters: $\gamma = 0.05$, $\Delta t = 0.05$, $N_{\text{step}} = 10$, and $K = 400$. The MBO scheme is allowed to run up to 100 iterations.