

SURF 2017 INTERIM REPORT 2

V. CHEN

ADVISED BY A. M. STUART AND M. M. DUNLOP

THIS SURF IS FUNDED BY THE AEROSPACE CORPORATION

1. Introduction. Over the past month, I have been implementing different hierarchical algorithms for the semi-supervised clustering problem. In the Bayesian formulation, recall the prior distribution,

$$(1) \quad u = \sum_{j=0}^M (\lambda_j + \tau^2)^{-\alpha/2} \xi_j q_j, \quad \xi_j \sim \mathcal{N}(0, 1) \quad \text{i.i.d.}$$

as well as the general form of the prior,

$$(2) \quad u = \sum_{j=0}^M f(\lambda_j) \xi_j q_j, \quad \xi_j \sim \mathcal{N}(0, 1) \quad \text{i.i.d.}$$

The models to be considered for this project evolve as follows:

For models (A), (B), and (C), the hyperparameters τ, α, M refer to the prior in (1)

(A) Use fixed $\tau, \alpha, M = N - 1$ to cluster data.

(B) Learn τ, α ; fix $M = N - 1$.

(C) Learn τ, α, M .

For models (D) and (E), samples from a new prior are given by $u = \sum_{j=0}^M v_j \xi_j q_j$.

(D) Learn $\{v_j\}_{j=0}^M$ with M fixed.

(E) Learn $\{v_j\}_{j=0}^M$ and M .

Model (F) can apply to all of the previous models.

(F) Multiclass, hierarchical on number of classes.

Model (G) focuses on learning the arbitrary function f defined in (2).

(G) Learn $f(\lambda), M$.

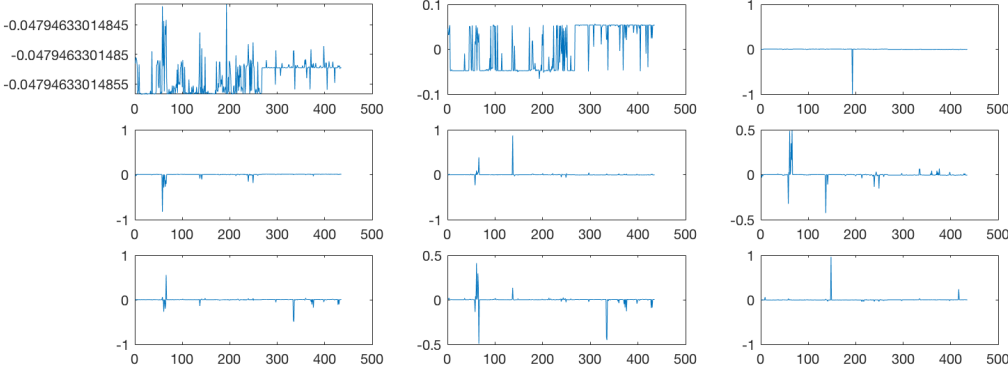
So far, I have implemented models (A) through (E) and have started working on (F) applied to (A) and (B). The past month has been devoted to experiments testing individual models and comparing their performances, using metrics such as classification accuracy and variance. Classification accuracy is well defined for the model problems that I tested the data on. These problems are outlined in [section 2](#).

2. Model problems. Applications to be considered include the following:

2.1. Voting records. A data set of the voting records on 16 votes of 435 individuals from the U.S. House of Representatives. The data is ordered so that all the representatives from the party are contiguous in the indexing. The graph has $N = 435$ nodes, with feature vectors that are $d = 16$ dimensional. Each component of the feature vectors is either $+1, -1$, or 0 to indicate voting for, voting against, or abstaining on a particular bill, respectively. On this data set, we want to perform binary classification with a subset of labeled data. For this problem, we will use the unnormalized Laplacian with fixed length-scale weights. That is, we will define $w_{ij} = \exp\left(-\frac{d(x_i, x_j)}{2l^2}\right)$, $d_{p,q}(x, y) = \|x - y\|_p^q$ with $p = 2, q = 2, l = 1$. [Figure 1](#) is part of the spectrum of this choice of Laplacian.

2.2. Two moons. This is a synthetic data set constructed with two half circles in \mathbb{R}^2 with radius one. These are embedded in \mathbb{R}^{100} , and data points are sampled from the circles with Gaussian noise distributed as $\mathcal{N}(0, \sigma^2)$ added to each of the 100 dimensions. In this data set, we will be generating realizations of two moons with $N = 1000$ or $N = 2000$ nodes, each associated with $d = 100$ dimensional feature vectors. This data is again ordered, with the first $N/2$ nodes generated from the first half circle, and the latter $N/2$ from the second. Again, here we want to

FIG. 1. *Lowest 9 eigenvectors of unnormalized L of voting records.*



perform binary classification into the two half circles from which the data was generated. We will be using the self-tuning, symmetric Laplacian for this data set, introduced in [5]. This Laplacian matrix has weights that infer the local spatial scale from the data, removing the need to choose a fixed length-scale parameter. The use of this self-tuning Laplacian in the two-moons data set seems to encode more information in the eigenvectors. Compare Figure 2 with Figure 3.

FIG. 2. *Lowest 9 eigenvectors of self-tuning symmetric L of two moons, $\sigma = 0.04$.*

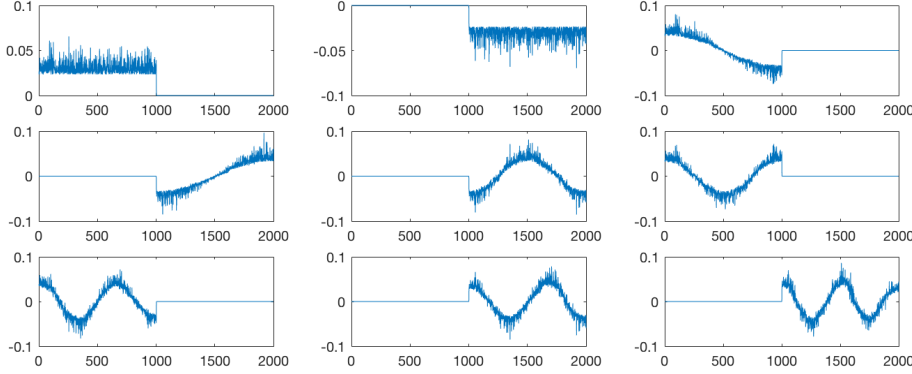
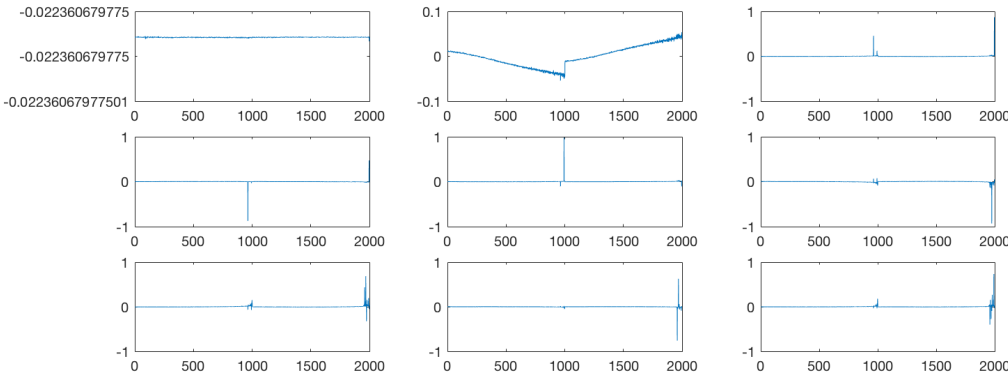
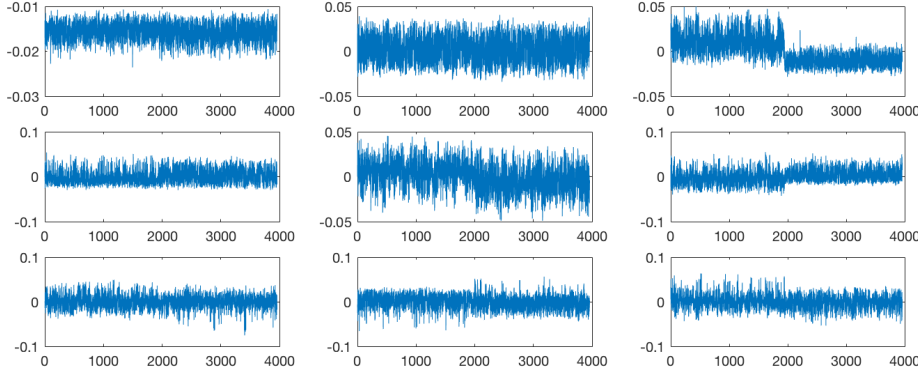


FIG. 3. *Lowest 9 eigenvectors of unnormalized fixed length L of two moons, $\sigma = 0.04$.*



2.3. MNIST data sets. This data set contains 70,000 images of 28×28 pixels with hand-written digits 0 through 9. The feature vectors are $d = 400$ dimensional and are formed by projection onto the first 50 PCA components. We will focus on binary classification between the digits 4 and 9. We will use the self-tuning, symmetric normalized Laplacian again. See Figure 4.

FIG. 4. Lowest 9 eigenvectors of self-tuning symmetric L of MNIST, digits 4 and 9.



3. Algorithms. This section discusses the algorithms that I have implemented for the hierarchical models.

3.1. Model (A). For the nonhierarchical model (A), I implemented a pCN algorithm first introduced in [2] and described in Algorithm 1.

Algorithm 1 General pCN adapted from [3]

- 1: Select $u^{(0)}$. Select τ, α . Select $\beta \in [0, 1]$
 - 2: **for** $k = 0$ to S **do**
 - 3: Sample v from the prior distribution given in (1)
 - 4: Set $\hat{u}^{(k)} = (1 - \beta^2)^{1/2} u^{(k)} + \beta v$
 - 5: Set $\alpha(u^{(k)} \rightarrow \hat{u}^{(k)}) = \min(1, \exp(\Phi(u^{(k)}) - \Phi(\hat{u}^{(k)})))$
 - 6: Set $u^{(k+1)} = \hat{u}^{(k)}$ with probability $\alpha(u^{(k)} \rightarrow \hat{u}^{(k)})$, and set $u^{(k+1)} = u^{(k)}$ otherwise
 - 7: **end for**
 - 8: **return** $\{u^{(k)}\}$
-

3.2. Model (B). To implement model (B), which is hierarchical, we need to sample conditional distributions that govern u, τ , and α . When we attempt to sample the posterior distribution in the hierarchical methods, we could use the Gibbs sampler. The basic form of Gibbs sampling has two repeated steps:

- Update $u^{(n+1)} \sim u | \theta^{(n)}, y$
- Update $\theta^{(n+1)} \sim \theta | u^{(n+1)}, y$

We use θ to represent the hyperparameters, which could include τ, α , and M . However, we cannot sample the conditional distributions directly, so we replace direct sampling with Markov Chain Monte Carlo (MCMC) indirect sampling methods, which are invariant with respect to each of the exact conditional distributions. With Metropolis-within-Gibbs, at every step we update $u^{(n+1)}, \tau^{(n+1)}$ and $\alpha^{(n+1)}$ each with MCMC to target these conditional distributions. One method is to update u, τ, α in a block, by proposing $(u, \theta) \rightarrow (\hat{u}, \hat{\theta})$ and computing the transition probability for this step.

The algorithms that we will implement will be slightly different, as we will independently propose $u \rightarrow \hat{u}, \tau \rightarrow \hat{\tau}, \alpha \rightarrow \hat{\alpha}, M \rightarrow \hat{M}$, and perform these transitions in separate steps. At each step, we fix all but the proposed parameter, and we compute the transition probability. This is the algorithm that we will be using for the following hierarchical algorithms.

Our first of two hierarchical algorithms for model (B) is deemed “centered,” as contrasted with the second algorithm described later in this section. Using Metropolis-within-Gibbs, we will require an expression for the posterior, so define $f(u, \tau, \alpha)$ as the joint posterior distribution. By Bayes’ theorem,

$$f(u, \tau, \alpha) \propto \exp(-\Phi(u)) \times \mathbb{P}(u, \tau, \alpha) = \exp(-\Phi(u)) \mathbb{P}(u|\theta) \pi_0(\theta).$$

Recall from (1) that the prior is distributed as $\mathbf{N}(0, C)$ with the covariance matrix $C(\theta) = C(\tau, \alpha) = (L + \tau^2 I)^{-\alpha}$. We can write

$$(3) \quad f(u, \tau, \alpha) \propto \frac{1}{\sqrt{(2\pi)^d \det C(\tau, \alpha)}} \exp \left(-\Phi(u) - \frac{1}{2} \langle u, C(\tau, \alpha)^{-1} u \rangle + \log(\pi_0(\tau, \alpha)) \right).$$

The normalization constant $\det(C(\theta))$ depends on τ, α now and does not cancel out. Using this expression for the posterior, I implemented the Metropolis-within-Gibbs method for model (B) in Algorithm 2.

Algorithm 2 Hierarchical on τ, α

- 1: Initialize $u^{(0)} = q_1$, the Fiedler vector expressed in the standard basis.
- 2: Initialize $\tau^{(0)}, \alpha^{(0)}$. Select $\beta \in [0, 1]$
- 3: Pick ϵ_1, ϵ_2 , the jump sizes for τ, α respectively.
- 4: **for** $k = 0$ to S **do**
- 5: Sample v from the prior distribution and expressed in the eigenbasis $\triangleright u|y, \tau, \alpha$.
- 6: Expressing u in the eigenbasis, set a proposal $\hat{u}^{(k)} = (1 - \beta^2)^{1/2} u^{(k)} + \beta v$
- 7: Set $u^{(k+1)} = \hat{u}^{(k)}$ with probability

$$A(u^{(k)} \rightarrow \hat{u}^{(k)}) = \min \left\{ 1, \exp(\Phi(u^{(k)}) - \Phi(\hat{u}^{(k)})) \right\}$$

- 8: Set a proposal $\hat{\tau}^{(k)} = \tau^{(k)} + \epsilon_1 t$ for $t \sim \mathbf{N}(0, 1)$ $\triangleright \tau|y, u, \alpha$
- 9: Set $\tau^{(k+1)} = \hat{\tau}^{(k)}$ with probability

$$A(\tau^{(k)} \rightarrow \hat{\tau}^{(k)}) = \min \left\{ 1, \frac{f(u^{(k+1)}, \hat{\tau}^{(k)}, \alpha^{(k)})}{f(u^{(k+1)}, \tau^{(k)}, \alpha^{(k)})} \right\}$$

(Using the eigenbasis representation of u for computation) $\triangleright f$ defined in (3)

- 10: Set a proposal $\hat{\alpha}^{(k)} = \alpha^{(k)} + \epsilon_2 a$ for $a \sim \mathbf{N}(0, 1)$ $\triangleright \alpha|y, u, \tau$
- 11: Set $\alpha^{(k+1)} = \hat{\alpha}^{(k)}$ with probability

$$A(\alpha^{(k)} \rightarrow \hat{\alpha}^{(k)}) = \min \left\{ 1, \frac{f(u^{(k+1)}, \tau^{(k+1)}, \hat{\alpha}^{(k)})}{f(u^{(k+1)}, \tau^{(k+1)}, \alpha^{(k)})} \right\}$$

- 12: **end for**
 - 13: **return** $\{u^{(k)}, \tau^{(k)}, \alpha^{(k)}\}$
-

We also looked at a different parameterization for model (B). Algorithm 3 uses the variable ξ that is related to the classifying function u by

$$(4) \quad T(\xi, \tau, \alpha) = \sum_{i=0}^M \frac{1}{(\lambda_i + \tau^2)^{\alpha/2}} \xi_i q_i = u$$

following (1). This algorithm is “non-centered” compared to the centered parameterization given in Algorithm 2, meaning that the unknown classifying variable ξ and the parameters $\theta = (\tau, \alpha)$ are a priori independent [4].

This model assumes the prior $\xi \sim \mathbf{N}(0, I)$. Similar to the centered case, define $g(\xi, \tau, \alpha)$ as the joint posterior distribution. By Bayes' theorem,

$$g(\xi, \tau, \alpha) \propto \exp(-\Phi(T(\xi, \tau, \alpha))) \times \mathbb{P}(\xi, \tau, \alpha) = \exp(-\Phi(T(\xi, \tau, \alpha))) \mathbb{P}(\xi) \pi_0(\tau, \alpha).$$

Note that $\mathbb{P}(\xi) = \frac{1}{\sqrt{(2\pi)^d \det I}} \exp(-\frac{1}{2}\langle \xi, I\xi \rangle) \propto \exp(-\frac{1}{2}\langle \xi, \xi \rangle)$ and the normalization constants drop out. We obtain:

$$(5) \quad g(\xi, \tau, \alpha) \propto \exp\left(-\Phi(T(\xi, \tau, \alpha)) - \frac{1}{2}\langle \xi, \xi \rangle + \log(\pi_0(\tau, \alpha))\right).$$

Algorithm 3 Non-centered parameterization: sampling ξ, τ, α

- 1: Choose $\xi^{(0)} \in \mathbb{R}^N, \alpha^{(0)} > 0, \beta \in (0, 1]$ and $\epsilon_1, \epsilon_2 > 0$.
- 2: **for** $k = 0$ to S **do**
- 3: Propose $\hat{\xi}^{(k)} = (1 - \beta^2)^{\frac{1}{2}} \xi^{(k)} + \beta \zeta^{(k)}, \zeta^{(k)} \sim \mathbf{N}(0, I)$
- 4: Make transition $\xi^{(k)} \rightarrow \hat{\xi}^{(k)}$ with probability

$$A(\xi^{(k)} \rightarrow \hat{\xi}^{(k)}) = \min \left\{ 1, \exp \left(\Phi(T(\xi^{(k)}, \tau^{(k)}, \alpha^{(k)})) - \Phi(T(\hat{\xi}^{(k)}, \tau^{(k)}, \alpha^{(k)})) \right) \right\}$$

$\triangleright T$ defined in (4)

- 5: Propose $\hat{\tau}^{(k)} = \tau^{(k)} + \epsilon_1 \rho^{(k)}, \rho^{(k)} \sim \mathbf{N}(0, 1)$
- 6: Make transition $\tau^{(k)} \rightarrow \hat{\tau}^{(k)}$ with probability

$$A(\tau^{(k)} \rightarrow \hat{\tau}^{(k)}) = \min \left\{ 1, \frac{g(\xi^{(k+1)}, \hat{\tau}^{(k)}, \alpha^{(k)})}{g(\xi^{(k+1)}, \tau^{(k)}, \alpha^{(k)})} \right\}$$

$\triangleright g$ defined in (5)

- 7: Propose $\hat{\alpha}^{(k)} = \alpha^{(k)} + \epsilon_2 \sigma^{(k)}, \sigma^{(k)} \sim \mathbf{N}(0, 1)$
- 8: Make transition $\alpha^{(k)} \rightarrow \hat{\alpha}^{(k)}$ with probability

$$A(\alpha^{(k)} \rightarrow \hat{\alpha}^{(k)}) = \min \left\{ 1, \frac{g(\xi^{(k+1)}, \tau^{(k+1)}, \hat{\alpha}^{(k)})}{g(\xi^{(k+1)}, \tau^{(k+1)}, \alpha^{(k)})} \right\}$$

9: **end for**

10: **return** $\{T(\xi^{(k)}, \tau^{(k)}, \alpha^{(k)}), \tau^{(k)}, \alpha^{(k)}\}$

3.3. Model (C). In this algorithm, we aim to learn M , the number of eigenvectors we want to use, as well. This algorithm seems promising given the problem we encountered with the irregularity of higher eigenvectors with the algorithms in model (B), which we resolved by truncating the number of eigenvectors. Essentially, we want to be hierarchical about our level of truncation. We take a uniform prior for M , typically $\mathbf{U}(1, 70)$ or similar. It is convenient to implement this algorithm with the non-centered approach since we do not need to worry about M affecting terms such as the determinant of the covariance. The new equation for T , which relates ξ, τ, α, M to u , is:

$$(6) \quad T(\xi, \tau, \alpha, M) = \sum_{i=0}^M \frac{1}{(\lambda_i + \tau^2)^{\alpha/2}} \xi_i q_i = u.$$

We can compute the joint posterior on ξ, τ, α, M similarly to the derivation of (5):

$$(7) \quad g(\xi, \tau, \alpha, M) \propto \exp\left(-\Phi(T(\xi, \tau, \alpha, M)) - \frac{1}{2}\langle \xi, \xi \rangle + \log(\pi_0(\tau, \alpha, M))\right).$$

The only modification from [Algorithm 3](#) is to add a random walk proposal for M and compute the acceptance probability as the ratio of the joint posteriors. This algorithm is given in [Algorithm 4](#).

Algorithm 4 Non-centered parameterization, hierarchical with τ, α, M

- 1: Choose $\xi^{(0)} \in \mathbb{R}^N, \tau^{(0)}, \alpha^{(0)}, M^{(0)} > 0, \beta \in (0, 1]$ and $\epsilon_1, \epsilon_2 > 0$.
- 2: **for** $k = 0$ to S **do**
- 3: Propose $\hat{\xi}^{(k)} = (1 - \beta^2)^{\frac{1}{2}} \xi^{(k)} + \beta \zeta^{(k)}, \zeta^{(k)} \sim \mathcal{N}(0, I)$
- 4: Make transition $\xi^{(k)} \rightarrow \hat{\xi}^{(k)}$ with probability

$$A(\xi^{(k)} \rightarrow \hat{\xi}^{(k)}) = \min \left\{ 1, \exp \left(\Phi(T(\xi^{(k)}, \tau^{(k)}, \alpha^{(k)}, M^{(k)})) - \Phi(T(\hat{\xi}^{(k)}, \tau^{(k)}, \alpha^{(k)}, M^{(k)})) \right) \right\}$$

$\triangleright T$ defined in (6)

- 5: Propose $\hat{\tau}^{(k)} = \tau^{(k)} + \epsilon_1 \rho^{(k)}, \rho^{(k)} \sim \mathcal{N}(0, 1)$
- 6: Make transition $\tau^{(k)} \rightarrow \hat{\tau}^{(k)}$ with probability

$$A(\tau^{(k)} \rightarrow \hat{\tau}^{(k)}) = \min \left\{ 1, \frac{g(\xi^{(k+1)}, \hat{\tau}^{(k)}, \alpha^{(k)}, M^{(k)})}{g(\xi^{(k+1)}, \tau^{(k)}, \alpha^{(k)}, M^{(k)})} \right\}$$

$\triangleright g$ defined in (7)

- 7: Propose $\hat{\alpha}^{(k)} = \alpha^{(k)} + \epsilon_2 \sigma^{(k)}, \sigma^{(k)} \sim \mathcal{N}(0, 1)$
- 8: Make transition $\alpha^{(k)} \rightarrow \hat{\alpha}^{(k)}$ with probability

$$A(\alpha^{(k)} \rightarrow \hat{\alpha}^{(k)}) = \min \left\{ 1, \frac{g(\xi^{(k+1)}, \tau^{(k+1)}, \hat{\alpha}^{(k)}, M^{(k)})}{g(\xi^{(k+1)}, \tau^{(k+1)}, \alpha^{(k)}, M^{(k)})} \right\}$$

- 9: Propose $\hat{M}^{(k)} = M^{(k)} + Q$, with jump Q distributed as $\mathbb{P}(Q = k) \propto \frac{1}{1+|k|}$, $|Q|$ bounded.

- 10: Make transition $M^{(k)} \rightarrow \hat{M}^{(k)}$ with probability

$$A(M^{(k)} \rightarrow \hat{M}^{(k)}) = \min \left\{ 1, \frac{g(\xi^{(k+1)}, \tau^{(k+1)}, \alpha^{(k+1)}, \hat{M}^{(k)})}{g(\xi^{(k+1)}, \tau^{(k+1)}, \alpha^{(k+1)}, M^{(k)})} \right\}$$

11: **end for**

12: **return** $\{T(\xi^{(k)}, \tau^{(k)}, \alpha^{(k)}, M^{(k)}), \tau^{(k)}, \alpha^{(k)}\}$

3.4. Model (D). This algorithm reparameterizes the problem in terms of the random vectors v and ξ . v_j modifies the scale of influence of q_j , the j th eigenvector of the graph Laplacian, on the classifying function u .

$$(8) \quad T(v, \xi) = \sum_{i=0}^M v_i \xi_i q_i = u.$$

Here, M is fixed. We take $\xi \sim \mathcal{N}(0, I)$. Using good estimates for τ, α , perhaps obtained by algorithms that learn τ, α , set the prior on v to be:

$$v_j \sim \mathcal{U} \left((1 - a)(\lambda_j + \tau^2)^{-\alpha/2}, (1 + a)(\lambda_j + \tau^2)^{-\alpha/2} \right)$$

where a is a fixed scalar.

Finally, we derive an expression for the posterior with Bayes' theorem.

$$\begin{aligned} \mathbb{P}(v, \xi | y) &\propto \mathbb{P}(y | v, \xi) \mathbb{P}(v, \xi) \\ &\propto \exp(-\Phi(T(v, \xi))) \mathbb{P}(v) \mathbb{P}(\xi) \\ &\propto \exp \left(-\Phi(T(v, \xi)) + \log(\pi_0(v)) - \frac{1}{2} \langle \xi, \xi \rangle \right). \end{aligned}$$

Let $h(v, \xi)$ denote the joint posterior on v and ξ . Then,

$$(9) \quad h(v, \xi) \propto \exp \left(-\Phi(T(v, \xi)) + \log(\pi_0(v)) - \frac{1}{2} \langle \xi, \xi \rangle \right).$$

The updates for this algorithm will be done as follows:

- Update $\xi^{(k+1)}|v^{(k)}, y$ with a pCN proposal.
- Update $v^{(k+1)}|\xi^{(k+1)}, y$ with a random walk on the uniform prior for v .

This algorithm is given in [Algorithm 5](#).

Algorithm 5 Non-centered parameterization, hierarchical with v

- 1: Choose $v^{(0)}, \xi^{(0)} \in \mathbb{R}^N, \beta \in (0, 1], \epsilon > 0$.
- 2: **for** $k = 0$ to S **do**
- 3: Propose $\hat{\xi}^{(k)} = (1 - \beta^2)^{\frac{1}{2}} \xi^{(k)} + \beta \zeta^{(k)}, \zeta^{(k)} \sim \mathbf{N}(0, I)$
- 4: Make transition $\xi^{(k)} \rightarrow \hat{\xi}^{(k)}$ with probability

$$A(\xi^{(k)} \rightarrow \hat{\xi}^{(k)}) = \min \left\{ 1, \exp \left(\Phi(T(v^{(k)}, \xi^{(k)})) - \Phi(T(v^{(k)}, \hat{\xi}^{(k)})) \right) \right\}$$

- 5: Propose $\hat{v}^{(k)} = v^{(k)} + \epsilon \rho^{(k)}, \rho^{(k)} \sim \mathbf{N}(0, I)$
- 6: **if** $\hat{v}_j^{(k)}$ is outside of its prior interval for any j **then**
- 7: Reject and set $v^{(k+1)} = v^{(k)}$.
- 8: **else**
- 9: Make transition $v^{(k)} \rightarrow \hat{v}^{(k)}$ with probability

$$\begin{aligned} A(v^{(k)} \rightarrow \hat{v}^{(k)}) &= \min \left\{ 1, \frac{h(\hat{v}^{(k)}, \xi^{(k+1)})}{h(v^{(k)}, \xi^{(k+1)})} \right\} \\ &= \min \left\{ 1, \exp \left(\Phi(T(v^{(k)}, \xi^{(k+1)})) - \Phi(T(\hat{v}^{(k)}, \xi^{(k+1)})) \right) \right\} \end{aligned}$$

- 10: **end if**
 - 11: **end for**
 - 12: **return** $\{T(v^{(k)}, \xi^{(k)}), v^{(k)}, \xi^{(k)}\}$
-

3.5. Model (E). This model extends model (D), trying to learn M as well. Define:

$$(10) \quad T(v, \xi, M) = \sum_{i=0}^M v_i \xi_i q_i = u.$$

Taking the uniform prior on v , we obtain:

$$(11) \quad h(v, \xi, M) \propto \exp \left(-\Phi(T(v, \xi, M)) + \log(\pi_0(v, M)) - \frac{1}{2} \langle \xi, \xi \rangle \right).$$

This approach is described in [Algorithm 6](#).

3.6. Multiclass. We describe the algorithm to solve the multiclass clustering problem. Recall that for $v \in \mathbb{R}^k$, we define $S(v) = e_p, \quad p = \operatorname{argmax}(v_r)$ where $\{e_1, \dots, e_k\}$ is the standard basis of \mathbb{R}^k . Now, the classifying function is of the form $u : Z \rightarrow \mathbb{R}^k$, or $u = (u^{[1]}, u^{[2]}, \dots, u^{[k]})$ with functions $u^{[j]} : Z \rightarrow \mathbb{R}$. Each of the $u^{[j]}$ has prior distribution given by $\mathbf{N}(0, C)$ with $C = (L + \tau^2 I)^{-\alpha}$. We retain

$$\Phi(u) = \frac{1}{2\gamma^2} \sum_{l \in Z'} \operatorname{norm}(y(l) - S(u(l)))^2.$$

We will implement the non-centered approach. Here, the variable is $\xi = (\xi^{[1]}, \dots, \xi^{[k]})$ where each $\xi^{[j]}$ is associated with $u^{[j]}$ by the following:

$$T(\xi^{[j]}) = \sum_{i=1}^M \frac{1}{(\lambda_j + \tau^2)^{\alpha/2}} \xi_i^{[j]} q_i = u^{[j]}.$$

Algorithm 6 Non-centered parameterization, hierarchical with v, M

- 1: Choose $v^{(0)}, \xi^{(0)} \in \mathbb{R}^N, M^{(0)}, \beta \in (0, 1], \epsilon > 0$.
- 2: **for** $k = 0$ to S **do**
- 3: Propose $\hat{\xi}^{(k)} = (1 - \beta^2)^{\frac{1}{2}} \xi^{(k)} + \beta \zeta^{(k)}, \zeta^{(k)} \sim \mathcal{N}(0, I)$
- 4: Make transition $\xi^{(k)} \rightarrow \hat{\xi}^{(k)}$ with probability

$$A(\xi^{(k)} \rightarrow \hat{\xi}^{(k)}) = \min \left\{ 1, \exp \left(\Phi(T(v^{(k)}, \xi^{(k)}, M^{(k)})) - \Phi(T(v^{(k)}, \hat{\xi}^{(k)}, M^{(k)})) \right) \right\}$$

- 5: Propose $\hat{v}^{(k)} = v^{(k)} + \epsilon \rho^{(k)}, \rho^{(k)} \sim \mathcal{N}(0, I)$
- 6: **if** $\hat{v}_j^{(k)}$ is outside of its prior interval for any j **then**
- 7: Reject and set $v^{(k+1)} = v^{(k)}$.
- 8: **else**
- 9: Make transition $v^{(k)} \rightarrow \hat{v}^{(k)}$ with probability

$$\begin{aligned} A(v^{(k)} \rightarrow \hat{v}^{(k)}) &= \min \left\{ 1, \frac{h(\hat{v}^{(k)}, \xi^{(k+1)}, M^{(k)})}{h(v^{(k)}, \xi^{(k+1)}, M^{(k)})} \right\} \\ &= \min \left\{ 1, \exp \left(\Phi(T(v^{(k)}, \xi^{(k+1)}, M^{(k)})) - \Phi(T(\hat{v}^{(k)}, \xi^{(k+1)}, M^{(k)})) \right) \right\} \end{aligned}$$

- 10: **end if**
- 11: Propose $\hat{M}^{(k)} = M^{(k)} + Q$, with jump Q distributed as $\mathbb{P}(Q = k) \propto \frac{1}{1+|k|}$, $|Q|$ bounded.
- 12: Make transition $M^{(k)} \rightarrow \hat{M}^{(k)}$ with probability

$$\begin{aligned} A(M^{(k)} \rightarrow \hat{M}^{(k)}) &= \min \left\{ 1, \frac{h(v^{(k+1)}, \xi^{(k+1)}, \hat{M}^{(k)})}{h(v^{(k+1)}, \xi^{(k+1)}, M^{(k)})} \right\} \\ &= \min \left\{ 1, \exp \left(\Phi(v^{(k+1)}, \xi^{(k+1)}, M^{(k)}) - \Phi(v^{(k+1)}, \xi^{(k+1)}, \hat{M}^{(k)}) \right) \right\} \end{aligned}$$

- 13: **end for**
 - 14: **return** $\{T(v^{(k)}, \xi^{(k)}), v^{(k)}, \xi^{(k)}\}$
-

Then, define:

$$(12) \quad T(\xi) = (T(\xi^{[1]}), \dots, T(\xi^{[k]})) = u.$$

Refer to [Algorithm 7](#).

We can extend this algorithm to be hierarchical with τ, α, M as in the binary case. Define T as

$$T(\xi^{[j]}, \tau, \alpha, M) = \sum_{i=1}^M \frac{1}{(\lambda_j + \tau^2)^{\alpha/2}} \xi_i^{[j]} q_i = u^{[j]}.$$

One choice we have to make is whether there should be one universal copy of τ, α, M for all of the k copies of the prior, or if each copy should have its own independently-evolving hyperparameters.

4. Experiments comparing hierarchical and nonhierarchical algorithms. Inspired by the experiment comparing the different clustering algorithms given in [1], we ran a similar set of experiments to determine if a hierarchical approach was beneficial to clustering. The clustering algorithms we want to compare are: Fiedler thresholding, nonhierarchical pCN, and two hierarchical algorithms for learning τ, α as well as the classifying function u . For each of the four algorithms tested, 50 random realizations of the two moons data set were generated with $r = 1, N = 2000, d = 100$, and varying σ . We used the same RNG seed for the trials across different algorithms so that these 50 realizations are consistent over the different algorithms tested. With each two moons realization, we tested the classification accuracy with varying

Algorithm 7 Multiclass, Metropolis-within-Gibbs updates

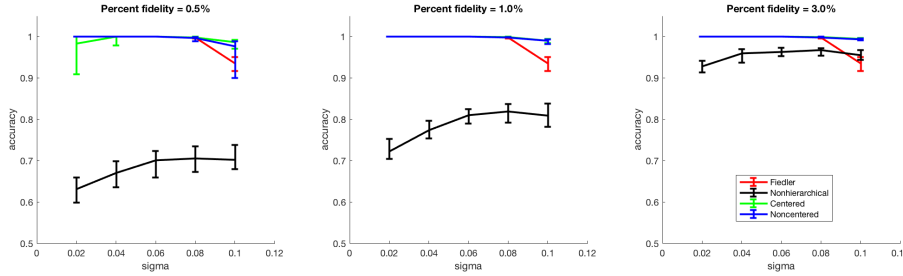
- 1: Choose $\xi^{(0)} = (\xi^{[1],(0)}, \dots, \xi^{[k],(0)}), \xi^{[j],(0)} \in \mathbb{R}^M$. Choose $\tau, \alpha, \beta \in (0, 1]$. \triangleright We refer to $\xi^{[j],(i)}$ as the i th observation of $\xi^{[j]}$.
- 2: **for** $i = 0$ to S **do**
- 3: **for** $j = 1$ to k **do**
- 4: Propose $\hat{\xi}^{[j],(i)} = (1 - \beta^2)^{\frac{1}{2}} \xi^{[j],(i)} + \beta \zeta^{(i)}, \zeta^{(i)} \sim \mathcal{N}(0, I)$
- 5: Define $\hat{\xi} = (\xi^{[1],(i+1)}, \dots, \xi^{[j-1],(i+1)}, \hat{\xi}^{[j],(i)}, \xi^{[j+1],(i)}, \dots)$
- 6: Make transition $\xi^{(i)} \rightarrow \hat{\xi}$ with probability

$$A(\xi^{(i)} \rightarrow \hat{\xi}) = \min \left\{ 1, \exp \left(\Phi(T(\xi^{(i)})) - \Phi(T(\hat{\xi})) \right) \right\}$$

- 7: **end for**
 - 8: **end for**
 - 9: **return** $\{T(\xi^{(i)})\}$
-

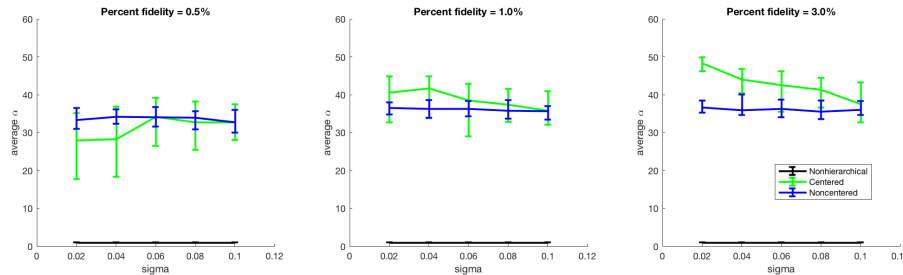
fidelity percents for the labeled data. Again, the chosen labeled set Z' is consistent across the trials for different algorithms. Our results are summarized in [Figure 5](#).

FIG. 5. Classification accuracy of different algorithms for two moons dataset compared with σ and percent fidelity. Plotted is the median classification accuracy with error bars that indicate the 25 and 75-th quantiles over the 50 trials for each parameter combination. For the pCN and hierarchical algorithms, we used spectral projection up to the first 50 eigenvectors. We ran 100000 iterations with a burn-in period of 1000, and we fixed $\gamma = 0.1$.



4.1. Learning τ and α . In these experiments, the hierarchical algorithms performed much better in classifying the two moons dataset than the nonhierarchical algorithm. We fixed $\tau = \alpha = 1$ in the nonhierarchical pCN, while we set $\tau^{(0)} = \alpha^{(0)} = 1$ for the hierarchical algorithms. In the hierarchical algorithms, we allowed $\tau \in [0.01, 60]$ and $\alpha \in [0.1, 60]$ as the uniform priors. The medians of the average values of α given by the hierarchical algorithms are shown in [Figure 6](#).

FIG. 6. Average values of α from the MCMC. Plotted is the median of the averages of α with error bars that indicate the 25 and 75-th quantiles over the 50 trials for each parameter combination. Notice that α is on average much larger in the hierarchical algorithms.



It appears that fixing $\alpha = 1$ in the nonhierarchical algorithm is partially responsible for the

discrepancy in performance. With small α , the prior on u drops off more slowly with the higher eigenvectors, while larger α enforces that these higher eigenvectors have less influence on u . This suggests that for the two moons data set, information for binary classification is concentrated in the lower eigenvectors, as expected.

5. Experiments with model (C).

5.1. M and σ relation. We tested [Algorithm 4](#) on the two moons dataset with varying σ , fixing 1% labeled nodes. We also fixed $\tau = 2, \alpha = 35$ by setting $\epsilon_1 = \epsilon_2 = 0$, so the algorithm is only learning ξ and M . We initialized $M^{(0)} = 50$, and the uniform prior of $M \sim \mathcal{U}(1, 70)$. We can see that for small $\sigma = 0.06$, M is small as only the first few eigenvectors are necessary (see [Figure 7](#)). However, when we choose a larger $\sigma = 0.2$, M needs to be larger (see [Figure 8](#)). For $\sigma = 0.06$, the classification accuracy was 100%. For $\sigma = 0.2$, the classification accuracy was 91.97%.

FIG. 7. $\sigma = 0.06$, trace of M

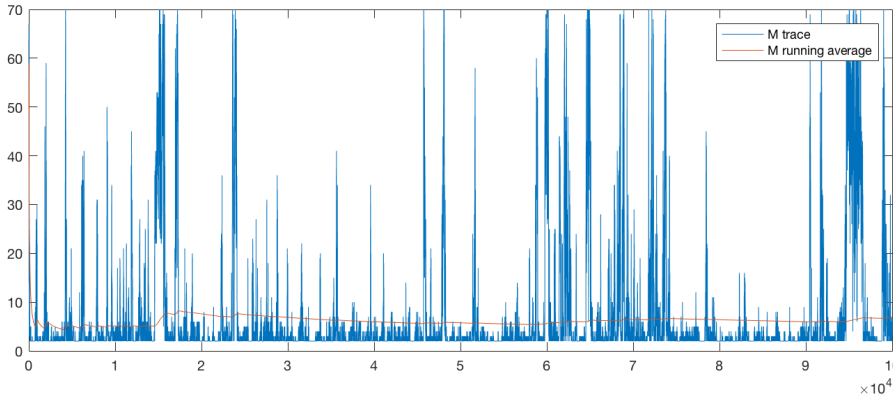
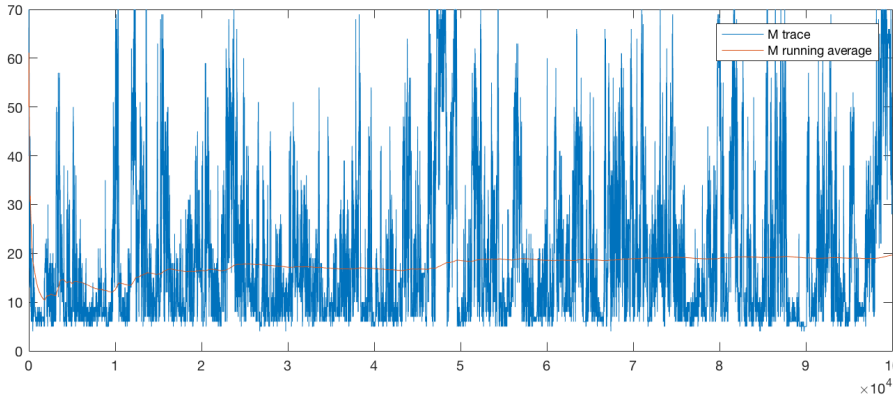


FIG. 8. $\sigma = 0.2$, trace of M



5.2. Nonhierarchical vs. Hierarchical. We will compare [Algorithm 1](#), pCN, with fixed τ and α against [Algorithm 4](#) with the same fixed τ and α , but learning M .

5.2.1. Voting records. The nonhierarchical algorithm for learning τ and α gave expected values of τ and α as $\tau \approx 2$ and $\alpha \approx 35$, so I chose those parameters for the pCN and for [Algorithm 4](#). For the pCN, all 435 eigenvectors were used. I narrowed the range of M allowed in [Algorithm 4](#) to 1 to 70. 5 labeled nodes were selected, consistent across the two methods. In this particular realization (seeded with `rng(5)`), the hierarchical algorithm achieves 87.74%

while the nonhierarchical algorithm achieves 87.67% classification accuracy. Figure 9, Figure 10, and Figure 11 show the results of the pCN. Figure 15, Figure 12, Figure 13 show the results the hierarchical algorithm. Notice in Figure 15 that there seems to be an important eigenvector for classification indexed around 30, as M seldom drops below 30. Looking at the eigenvectors around index 30, I plotted the 34th eigenvector in Figure 14. From a purely visual perspective, it does appear that this eigenvector corresponds decently to the final classification in Figure 12 by looking at where some of the “spikes” are. After trying some other realizations, it seems that being hierarchical may have some benefits, but more experiments are needed.

FIG. 9. *Algorithm 1*, average of $S(u)$

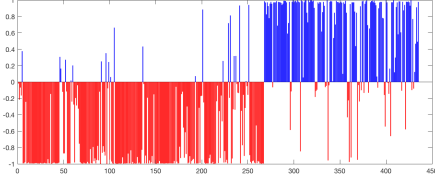


FIG. 10. *Algorithm 1*, u acceptance probability

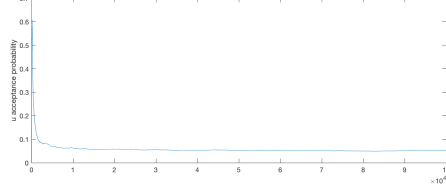


FIG. 11. *Algorithm 1*, running average of $S(u(i))$ for select i

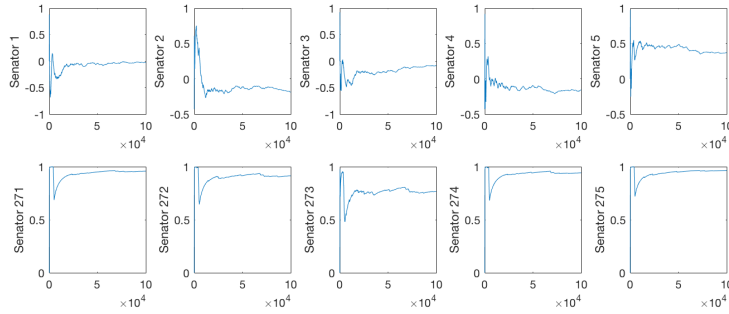


FIG. 12. *Algorithm 4*, average of $S(u)$

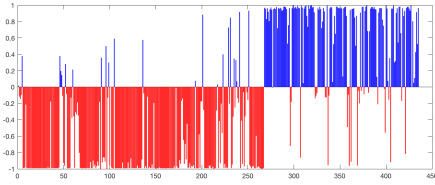


FIG. 13. *Algorithm 4*, M and ξ acceptance probabilities

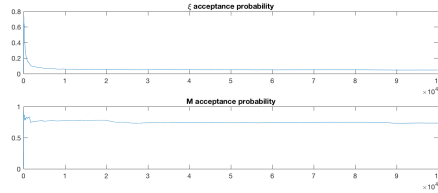


FIG. 14. 34th eigenvector of unnormalized fixed length-scale L

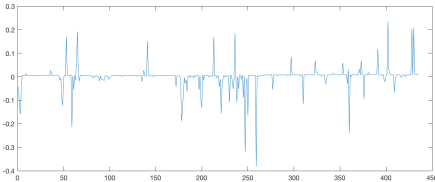
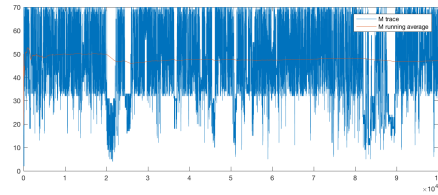


FIG. 15. *Algorithm 4*, trace of M



5.2.2. Two moons. I tried a similar set of experiments on the two moons data. I fixed $\tau = 2, \alpha = 35$, again following the expected values of these hyperparameters from the noncentered algorithm that learns τ, α . I set the pCN to use the first 100 eigenvectors, and I set the range for M in Algorithm 4 to be $[1, 70]$. The two moons data was generated with $N = 2000, d = 100, \sigma = 0.2$ and 1% fidelity. The nodes selected for fidelity are random but consistent between the different methods, as is the data itself. The colored diamonds in the scatter plots are the labeled nodes. One realization of this test is shown in Figure 16, Figure 17 for pCN, and Figure 18, Figure 19, Figure 20 for the hierarchical algorithm. In this particular realization, the two algorithms performed similarly, with 90.56% accuracy for pCN and 91.97% accuracy for the hierarchical algorithm as shown in the final clusterings in Figure 16 and Figure 18. In Figure 20, M seems to have moved from its uniform prior distribution, preferring a mean of around 20. The results again suggest that there is some value to being hierarchical with M even with good choices of τ and α , but more experiments are needed.

FIG. 16. *pCN, final classification projected into first two dimensions.*



FIG. 17. *pCN, u acceptance probability.*

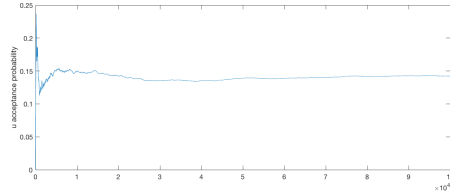


FIG. 18. *Hierarchical, final classification projected into first two dimensions.*

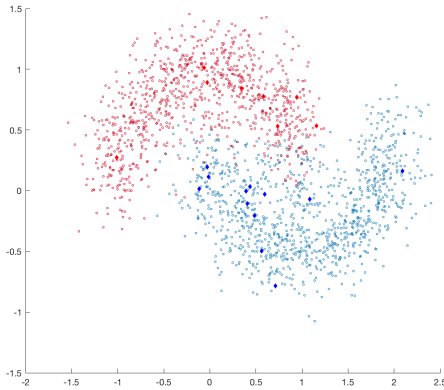
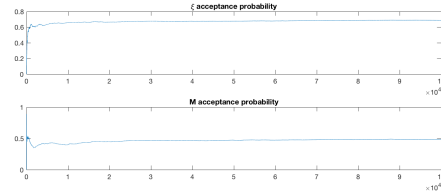


FIG. 19. *Hierarchical, ξ acceptance probability.*



6. Experiments with model (E). We run some experiments with Algorithm 6, model (E), on the two moons dataset. We again fix $\tau = 2, \alpha = 35$, and pick $a = 0.5$ to create the prior on v . We pick the walk size for v to be $\epsilon = 1.0 \times 10^{-13}$. We allow the prior on M to be uniform, $U(1, 70)$. If we compare this algorithm against model (C) with fixed τ, α , this is

FIG. 20. Hierarchical, seed $\text{rng}(3)$. M trace.

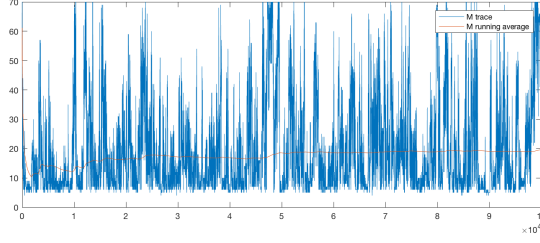
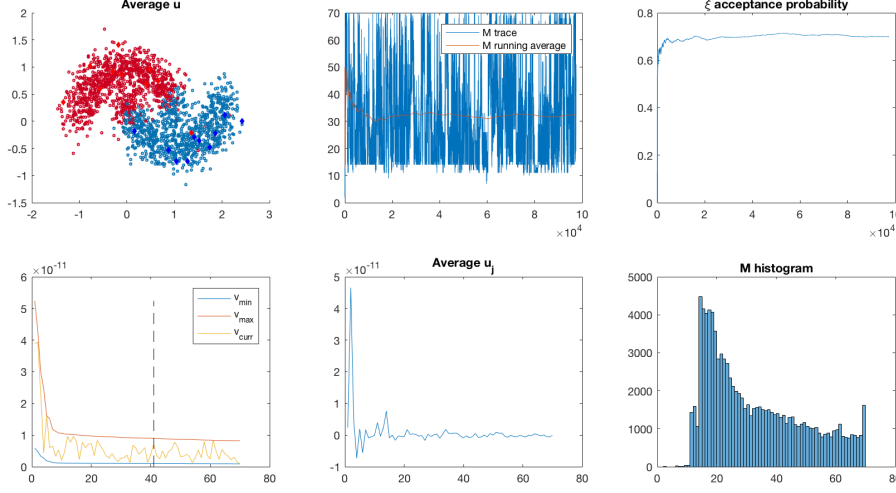


FIG. 21. Learning v, M . The bottom left figure was one observation of the v_j .



equivalent to setting $\epsilon = 0$ and fixing v to be at the mean of its prior. Learning v and M does not seem to greatly affect the classification accuracy compared to just learning M . Consider the following run on two moons, with fidelity = 1%, $\sigma = 0.2$, in Figure 21 and Figure 22. The model that learns v, M obtained 85.45% accuracy while the one that only learns M obtained 86.52% accuracy. The histograms of M appear very similar, with the same cutoff observed at the 14th index. The average values of the u_j coefficients also seem to agree in magnitude and sign between the two algorithms.

7. Multiclass experiments. We use the MNIST dataset and attempt to cluster k different digits. We want to compare the algorithm that learns τ, α against the nonhierarchical algorithm with fixed τ, α , so we fix $\tau = 2, \alpha = 35$ for the nonhierarchical algorithm and initialize $\tau = 2, \alpha = 35$ for the hierarchical algorithm. Additionally, for the hierarchical algorithm we currently restrict (τ, α) to be the same for each of the k priors. We fix $\gamma = 0.1, \beta = 0.1$ for both algorithms. We test on the digits [1, 4, 9] with the same 1% fidelity selected between the two trials. The final classification is given by $S(\mathbb{E}(S(u)))$, with a burn-in of 5000 and 100000 total iterations. The hierarchical algorithm obtained 94.33% accuracy while the nonhierarchical algorithm obtained 88.67%. This observation supports the claim that there is value in learning τ, α . See Figure 23, Figure 24, Figure 25, Figure 26, Figure 27.

Using $S(\mathbb{E}(S(u)))$ seems to give a better accuracy for the hierarchical algorithm than $S(\mathbb{E}(u))$, which would only obtain 88.62% accuracy for the above example. This seems to have to do with the high variance in the scale of u , which allows a few observations to dominate the classification. However, the choice of $S(\mathbb{E}(u))$ is more natural perhaps, since $\mathbb{E}(u)$ is the expected value of u in the posterior distribution.

FIG. 22. *Learning only M*

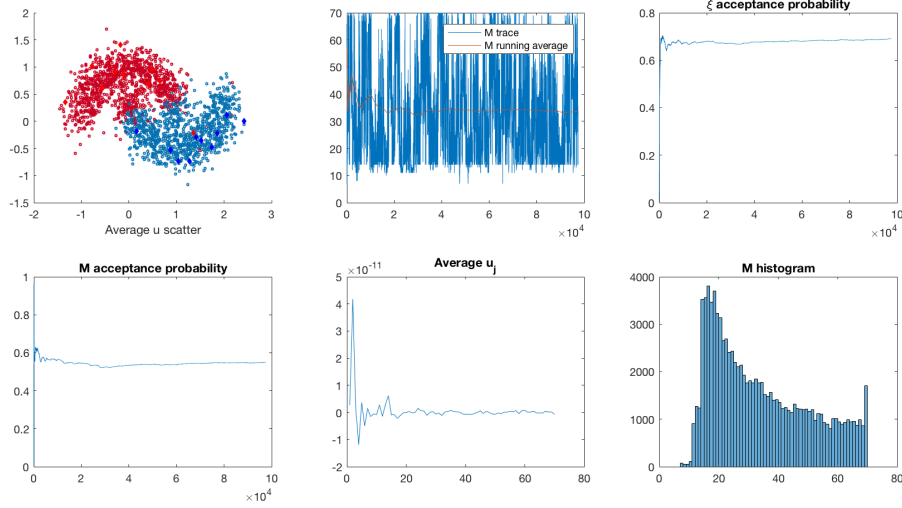


FIG. 23. *Nonhierarchical multiclass [1, 4, 9], final classification*

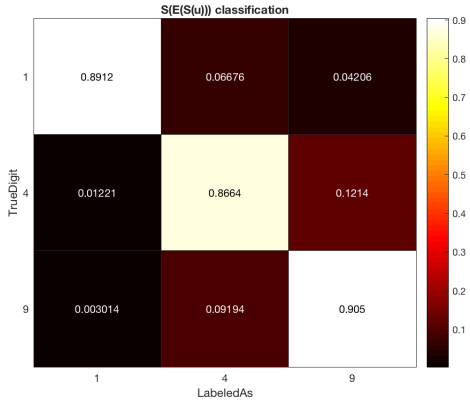


FIG. 24. *Nonhierarchical multiclass [1, 4, 9], ξ acceptance probability*

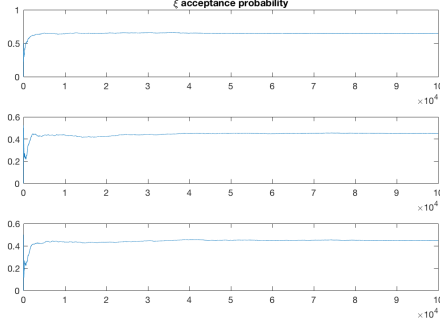


FIG. 25. *Hierarchical multiclass [1, 4, 9], final classification*

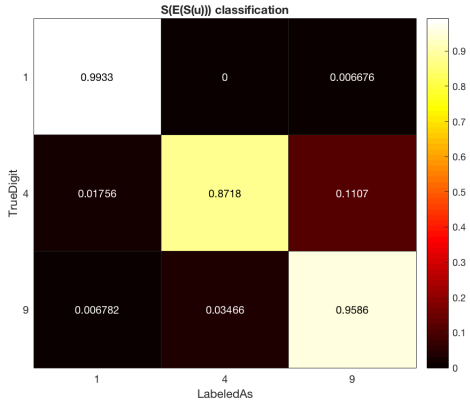


FIG. 26. *Hierarchical multiclass [1, 4, 9], ξ acceptance probability*

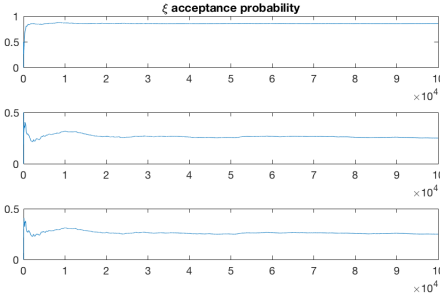


FIG. 27. Hierarchical multiclass [1, 4, 9], τ and α traces

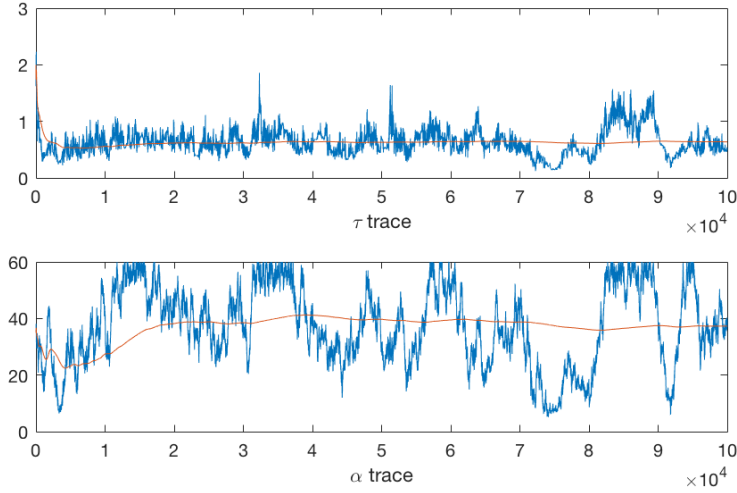
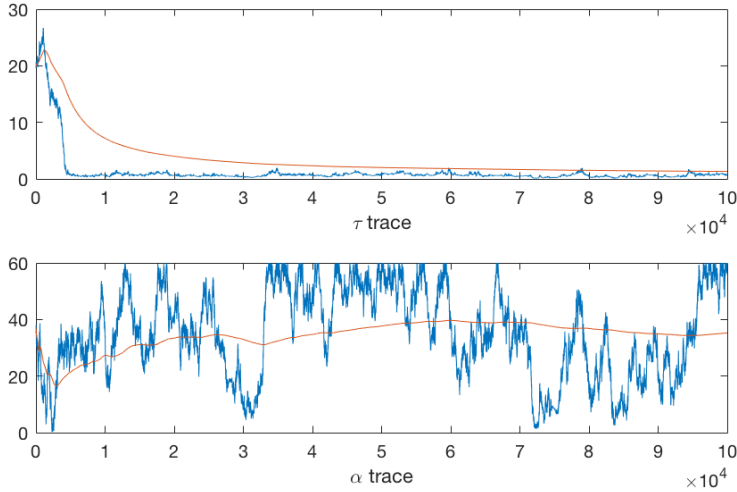


FIG. 28. Hierarchical multiclass [1, 4, 9], $\tau_0 = 20$ and learning τ



We can also see the same behavior with τ that we observed in the binary case. Initializing $\tau = 20$ and $\alpha = 35$ in [Figure 28](#), we can see that τ learns that it should be small, averaging around 0.6. This value seems pretty consistent across different initial values of τ . We can also see that there is strong correlation between τ and α , for instance by visually matching the extrema in [Figure 27](#).

8. Challenges. One problem I have been running into is that these MCMC simulations take a while to run. Since a large number of iterations is needed to ensure convergence and there is matrix multiplication for converting from eigenbasis to standard basis, large test suites are time consuming. My workaround for this problem is to run these tests overnight and analyze the results in the morning.

Another problem I encountered was the difference between $S(\mathbb{E}(S(u)))$ and $S(\mathbb{E}(u))$ as the choice of the classifying function. While $S(\mathbb{E}(u))$ is more natural as discussed in [section 7](#), this allows large u to dominate the entire classification, and large u can arise when τ is small, α is large, and the choice is accepted. I noticed this issue when I saw that the classification sometimes would

get stuck at a percentage with little change for the remaining iterations. My solution has been to use $S(\mathbb{E}(S(u)))$ for now, and it seems to give better classification.

9. Goals. One research goal for the remainder of the summer is to gather more data comparing the algorithm that learns M against the nonhierarchical algorithm to determine if the method is beneficial. It will also be interesting to better analyze the relationship between σ and M in the two moons dataset. Preliminary experiments suggest that in low σ data, small M is enough for an accurate clustering, but in large σ , higher M is necessary as more eigenvectors are needed to fit the problem. It could be interesting to create a plot of average M versus increasing σ to test this hypothesis.

Also, I hope to look at different prior choices for v for [Algorithm 5](#) and [Algorithm 6](#). It appears that learning v_j and M with a uniform prior on v_j is similar to just learning M , so a different prior choice could be useful. For example, one idea is to define $l_j = \frac{1}{v_j^2}$ and assume a gamma distribution for l_j with shape depending on j .

Another goal is to study different multiclass hierarchical methods and implement a multiclass algorithm that learns M . I am working on implementing the different choices described in [subsection 3.6](#). In particular, independently evolving hyperparameters τ, α, M for each of the k priors should be studied.

REFERENCES

- [1] A. L. BERTOZZI, X. LUO, A. M. STUART, AND K. C. ZYGALAKIS, *Uncertainty quantification in the classification of high dimensional data*, ArXiv e-prints, (2017), <https://arxiv.org/abs/1703.08816>.
- [2] A. BESKOS, G. O. ROBERTS, A. M. STUART, AND J. VOSS, *MCMC methods for diffusion bridges*, Stochastics and Dynamics, 8 (2008), pp. 319–350.
- [3] 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, 28 (2013), pp. 424–446.
- [4] O. PAPASPILIOPOULOS, G. O. ROBERTS, AND M. SKÖLD, *A general framework for the parametrization of hierarchical models*, Statistical Science, (2007), pp. 59–73.
- [5] L. ZELNIK-MANOR AND P. PERONA, *Self-tuning spectral clustering*, in Advances in neural information processing systems, 2005, pp. 1601–1608.