

A penalized likelihood method for classification with matrix-valued predictors

Aaron J. Molstad* and Adam J. Rothman
School of Statistics, University of Minnesota

Abstract

We propose a penalized likelihood method to fit the linear discriminant analysis model when the predictor is matrix valued. We simultaneously estimate the means and the precision matrix, which we assume has a Kronecker product decomposition. Our penalties encourage pairs of response category mean matrix estimators to have equal entries and also encourage zeros in the precision matrix estimator. To compute our estimators, we use a blockwise coordinate descent algorithm. To update the optimization variables corresponding to response category mean matrices, we use an alternating minimization algorithm that takes advantage of the Kronecker structure of the precision matrix. We show that our method can outperform relevant competitors in classification, even when our modeling assumptions are violated. We analyze an EEG dataset to demonstrate our method’s interpretability and classification accuracy.

Keywords: alternating minimization algorithm, classification, penalized likelihood

1 Introduction

We propose a method for classification when the predictor is matrix valued, e.g. classification of hand-written letters. Standard vector-valued predictor classification methods, such as logistic regression and linear discriminant analysis, could be applied, but they would not take advantage of the matrix structure.

Logistic regression based methods for classification with a matrix-valued predictor have been proposed. Zhou and Li (2014) proposed a nuclear norm penalized likelihood estimator of the regression coefficient matrix $B_* \in \mathbb{R}^{r \times c}$ in a generalized linear model, where the value of the matrix predictor $x \in \mathbb{R}^{r \times c}$ enters the model through the trace of $B_*^T x$. In the same setup, Hung and Wang (2013) assumed that $\text{vec}(B_*) = \beta_* \otimes \alpha_*$ where vec stacks the columns of its argument, $\alpha_* \in \mathbb{R}^r$, $\beta_* \in \mathbb{R}^c$, and \otimes is the Kronecker product. This decomposition was also studied in the dimension reduction literature (Li et al., 2010).

There also exist non-likelihood based methods for classification with a matrix-valued predictor. These approaches modify Fisher’s linear discriminant criterion, e.g. 2D-LDA (Li and Yuan, 2005), matrix discriminant analysis (Zhong and Suslick, 2015), and penalized matrix discriminant analysis (Zhong and Suslick, 2015).

We propose a penalized likelihood method for classification with a matrix-valued predictor. Our method estimates the parameters in the linear discriminant analysis model. Let $x_i \in \mathbb{R}^{r \times c}$ be the

*Corresponding author: molst029@umn.edu

measured predictor for the i th subject and let $y_i \in \{1, \dots, J\}$ be the measured categorical response for the i th subject ($i = 1, \dots, n$). We assume that $(x_1, y_1), \dots, (x_n, y_n)$ are a realization of n independent copies of (X, Y) with the following distribution. The marginal distribution of Y is defined by $P(Y = j) = \pi_j$ ($j = 1, \dots, J$), where the π_j 's are unknown; and

$$\text{vec}(X) \mid Y = j \sim N_{rc} \{ \text{vec}(\mu_{*j}), \Sigma_* \}, \quad j = 1, \dots, J, \quad (1)$$

where $\mu_{*j} \in \mathbb{R}^{r \times c}$ is the unknown mean matrix for the j th response category, and Σ_* is the unknown rc by rc covariance matrix.

We make the simplifying assumption that

$$\Sigma_*^{-1} = \Delta_* \otimes \Phi_*, \quad (2)$$

which is equivalent to $\Sigma_* = \Delta_*^{-1} \otimes \Phi_*^{-1}$, where Φ_* is an unknown r by r precision matrix with $\sum_{a,b} |\Phi_{*a,b}| = r$, and Δ_* is an unknown c by c precision matrix. The norm condition on Φ_* is added for identifiability: see Roś et al. (2016) for more on identifiability under (2). This simplification of a covariance matrix makes the conditional distributions in (1) become matrix normal (Gupta and Nagar, 2000). This exploits the matrix structure of the predictor by reducing the number of parameters in the precision matrix from $O(r^2 c^2)$ to $O(r^2 + c^2)$.

Several authors have proposed and studied penalized likelihood estimators of Φ_* and Δ_* when $J = 1$ (Allen and Tibshirani, 2010; Zhang and Schneider, 2010; Tsiligkaridis et al., 2012; Leng and Tang, 2012; Zhou, 2014).

In this paper, we propose a penalized likelihood method to fit (1) with the assumption in (2). Our penalties encourage fitted models that can be easily interpreted by practitioners. We use a blockwise coordinate descent algorithm to compute our estimators. To exploit (2) computationally, we use an alternating minimization algorithm (Tseng, 1991) in one of our block updates. This algorithm scales more efficiently than other popular algorithms, which makes our method computationally feasible for high-dimensional problems. We show that our algorithm has the same computational complexity order as the unpenalized likelihood version, which also requires a blockwise coordinate descent algorithm (Dutilleul, 1999).

2 Penalized likelihood estimation

2.1 Proposed method

Let \mathbb{S}_+^m be the set of symmetric and positive definite m by m matrices. The maximum likelihood estimators of the μ_{*j} 's, Φ_* , and Δ_* minimize the function $g : (\mathbb{R}^{r \times c})^J \times \mathbb{S}_+^r \times \mathbb{S}_+^c \rightarrow \mathbb{R}$ defined by

$$g(\mu, \Phi, \Delta) = \frac{1}{n} \sum_{j=1}^J \left[\sum_{i=1}^n 1(y_i = j) \text{tr} \{ \Phi(x_i - \mu_j) \Delta (x_i - \mu_j)^T \} \right] - c \log \det(\Phi) - r \log \det(\Delta),$$

where $\mu = (\mu_1, \dots, \mu_J)$. We propose the penalized likelihood estimators defined by

$$\begin{aligned} (\hat{\mu}, \hat{\Delta}, \hat{\Phi}) = \arg \min_{(\mu, \Phi, \Delta) \in \mathcal{T}} & \left\{ g(\mu, \Phi, \Delta) + \lambda_1 \sum_{j < m} \|w_{j,m} \circ (\mu_j - \mu_m)\|_1 + \lambda_2 \|\Delta \otimes \Phi\|_1 \right\}, \quad (3) \\ & \text{subject to } \|\Phi\|_1 = r \end{aligned}$$

where $\mathcal{T} = (\mathbb{R}^{r \times c})^J \times \mathbb{S}_+^r \times \mathbb{S}_+^c$; \circ is the Hadamard product; $\|\cdot\|_1$ is the sum of the absolute values of the entries of its argument; λ_1 and λ_2 are nonnegative tuning parameters; and the $w_{j,m}$'s are r by c user-specified weight matrices.

The first penalty in (3) encourages solutions for which pairs of the mean matrix estimates have some equal entries, where this equality occurs in the same locations. Without the first penalty, i.e. $\lambda_1 = 0$, the proposed estimators of the μ_{*j} 's are sample mean matrices. If $\lambda_1 > 0$, then the proposed estimators of the μ_{*j} 's are affected by the estimators of Φ_* and Δ_* .

We recommend selecting weights similar to those prescribed by Guo (2010). We suggest using

$$w_{j,m}^{-1} = |\bar{x}_j - \bar{x}_m|, \quad 1 \leq j < m \leq J$$

where $\bar{x}_j = \sum_{i=1}^n 1(y_i = j)x_i$. Alternatively, one could use weights based on t -test statistics or could use weights that incorporate prior information.

The second penalty in (3) has a simple impact: for sufficiently large values of λ_2 , some of the entries in the estimate of $\Delta_* \otimes \Phi_*$ are zero, which occurs if and only if either the estimate of Δ_* or the estimate of Φ_* has some zero entries. To encourage zeros in estimates of Φ_* or Δ_* separately, one could use two separate L_1 penalties. Our computational algorithm can be easily adapted to accommodate this case.

The tuning parameters λ_1 and λ_2 can be chosen by minimizing the misclassification rate on a validation set.

2.2 Related work

Xu et al. (2015) proposed fitting the standard linear discriminant analysis model for a vector-valued predictor by penalized likelihood. We can express their parameter estimates in our matrix-predictor setup by setting the number of columns of the matrix predictor to one. Specifically, with $c = 1$ and $\Delta = 1$, Xu et al. (2015) parameter estimates are

$$\arg \min_{(\mu, \Phi) \in (\mathbb{R}^r)^J \times \mathbb{S}_+^r} \left\{ g(\mu, \Phi, 1) + \lambda_1 \sum_{j < m} \|w_{j,m} \circ (\mu_j - \mu_m)\|_1 + \lambda_2 \sum_{a \neq b} |\Phi_{ab}| \right\}. \quad (4)$$

One could view our method as the matrix-valued predictor extension of the method of Xu et al. (2015). Guo (2010) proposed a method that solves a restricted version of (4), where Φ is fixed at a diagonal matrix with pooled sample precision estimates on its diagonal.

Computationally, the algorithms proposed by Xu et al. (2015) and Guo (2010) for solving (4) suffer from numerical instability and do not scale efficiently for application to (3). In our simulation studies, we compare our proposed method to several competitors, including the method of Guo (2010). The method of Xu et al. (2015) is too slow computationally for the dimensions we consider,

so we only use it in a special case when Σ_* is known.

3 Computation

3.1 Overview

To solve (3), we use a block-wise coordinate descent algorithm. Each block update is a convex optimization problem. In the subsequent subsections, we show that updates for Φ and Δ can be expressed as the well-studied L_1 -penalized normal likelihood precision matrix estimation problem. We also use an alternating minimization algorithm for the block update for μ . The algorithm to compute our estimator, along with a set of auxiliary functions, is available in the R package `MatrixLDA`, which is included in the supplemental material.

3.2 Updates for Φ and Δ

We first derive the update for Φ . Define $\text{GL}(S, \tau)$ as

$$\text{GL}(S, \tau) = \arg \min_{\Theta \in \mathbb{S}_+} \{ \text{tr}(S\Theta) - \log |\Theta| + \tau \|\Theta\|_1 \}, \quad (5)$$

where S is some given nonnegative definite matrix and τ is a nonnegative tuning parameter. The optimization problem in (5) is the L_1 -penalized normal likelihood precision matrix estimation problem. Many algorithms and efficient software exist to solve (5): one good example is the graphical-lasso of Friedman et al. (2008).

Let f be the objective function in (3). Suppose Δ and μ are fixed. The minimizer of f with respect to Φ is

$$\tilde{\Phi} = \arg \min_{\Phi \in \mathbb{S}_+^r} \frac{1}{n} \sum_{j=1}^J \left[\sum_{i=1}^n 1(y_i = j) \text{tr} \{ \Phi(x_i - \mu_j) \Delta(x_i - \mu_j)^T \} \right] - c \log \det(\Phi) + \lambda_2 \|\Phi \otimes \Delta\|_1. \quad (6)$$

Using the fact that $\|\Phi \otimes \Delta\|_1 = \|\Phi\|_1 \|\Delta\|_1$ and

$$\frac{1}{n} \sum_{j=1}^J \left[\sum_{i=1}^n 1(y_i = j) \text{tr} \{ \Phi(x_i - \mu_j) \Delta(x_i - \mu_j)^T \} \right] = c \text{tr} \{ \Phi S_\phi(\mu, \Delta) \},$$

where

$$S_\phi(\mu, \Delta) = \frac{1}{nc} \sum_{j=1}^J \left\{ \sum_{i=1}^n 1(y_i = j) (x_i - \mu_j) \Delta (x_i - \mu_j)^T \right\},$$

we can express (6) as

$$\arg \min_{\Phi \in \mathbb{S}_+^r} \left[\text{tr} \{ \Phi S_\phi(\mu, \Delta) \} - \log |\Phi| + \frac{\lambda_1 \|\Delta\|_1}{c} \|\Phi\|_1 \right] = \text{GL} \left\{ S_\phi(\mu, \Delta), \frac{\lambda_1 \|\Delta\|_1}{c} \right\}.$$

After computing $\tilde{\Phi}$ with Δ fixed, we can enforce the constraint $\|\Phi\|_1 = r$ using a simple normal-

ization: we replace $(\tilde{\Phi}, \Delta)$ with $(\bar{\Phi}, \bar{\Delta})$, where

$$\bar{\Phi} = \frac{r}{\|\tilde{\Phi}\|_1} \tilde{\Phi}, \quad \bar{\Delta} = \frac{\|\tilde{\Phi}\|_1}{r} \Delta.$$

This ensures that $\|\bar{\Phi}\| = r$ without changing the objective function because $f(\mu, \Delta, \tilde{\Phi}) = f(\mu, \bar{\Delta}, \bar{\Phi})$.

Using a similar argument, the minimizer of f with respect to Δ with μ and Φ fixed is

$$\tilde{\Delta} = \text{GL} \left\{ S_\delta(\mu, \Phi), \frac{\lambda_1 \|\Phi\|_1}{r} \right\},$$

where

$$S_\delta(\mu, \Phi) = \frac{1}{nr} \sum_{j=1}^J \left\{ \sum_{i=1}^n 1(y_i = j) (x_i - \mu_j)^T \Phi (x_i - \mu_j) \right\}.$$

3.3 Update for μ

Let Δ and Φ be fixed. The minimizer of f with respect to μ is

$$\arg \min_{\mu \in \mathbb{R}^{(r \times c)J}} \frac{1}{n} \sum_{j=1}^J \left\{ \sum_{i=1}^n 1(y_i = j) \text{tr} [\Phi (x_i - \mu_j) \Delta (x_i - \mu_j)^T] \right\} + \lambda_1 \sum_{j < m} \|w_{j,m} \circ (\mu_j - \mu_m)\|_1. \quad (7)$$

Special cases of (7) have been solved using a majorize-minimize (MM) algorithm, where the penalty is majorized by its local-quadratic approximation at the current iterate (Hunter and Li, 2005). For example, Xu et al. (2015) solved (7) when $c = 1$ and $\Delta = 1$; and Guo (2010) solved (7) when $c = 1$, $\Delta = 1$, and Φ was diagonal. However, this MM algorithm suffers from numerical instability when iterates for μ_j and μ_m are similar from some (j, m) . Moreover, if we were to apply the MM algorithm to solve (7), then each iteration would have worst case computational complexity $O(r^2 c^2)$.

Instead of using an MM algorithm, we use an alternating minimization algorithm (Tseng, 1991; Chi and Lange, 2015) to solve (7). Our algorithm for solving (7) is more numerical stable, each iteration has worst case computational complexity $O(r^2 c + c^2 r)$ when distributed over $\max\{J, J(J-1)/2\}$ machines, and has a quadratic rate of convergence when implemented with the accelerations proposed by Goldstein et al. (2014). Both the MM algorithm and our alternating minimization algorithm require one eigendecomposition of Δ and of Φ .

Similarly to the setup of the ADMM algorithm (Boyd et al., 2011), we first express (7) as a constrained optimization problem:

$$\begin{aligned} & \underset{(\mu, \Theta) \in \mathcal{G}}{\text{minimize}} \quad g(\mu, \Phi, \Delta) + \lambda_1 \sum_{j < m} \|w_{j,m} \circ \Theta_{j,m}\| \\ & \text{subject to} \quad \Theta_{j,m} = \mu_j - \mu_m \quad 1 \leq j < m \leq J, \end{aligned} \quad (8)$$

where $\mathcal{G} = \mathbb{R}^{(r \times c)J} \times \mathbb{R}^{(r \times c)J(J-1)/2}$ and $\Theta = (\Theta_{1,2}, \dots, \Theta_{J-1,J})$. The augmented Lagrangian

for (8), using notation similar to Chi and Lange (2015), is

$$\begin{aligned}\mathcal{F}_\rho(\mu, \Theta, \Gamma) = & g(\mu, \Phi, \Delta) + \lambda_1 \sum_{j < m} \|w_{j,m} \circ \Theta_{j,m}\|_1 \\ & + \sum_{j < m} \text{tr} \{ \Gamma_{j,m}^T (\Theta_{j,m} - \mu_j + \mu_m) \} + \frac{\rho}{2} \sum_{j < m} \|\Theta_{j,m} - \mu_j + \mu_m\|_F^2,\end{aligned}$$

for step size parameter $\rho > 0$ and Lagrangian variables $\Gamma_{j,m} \in \mathbb{R}^{r \times c}$ for $1 \leq j < m \leq J$. Letting the superscript t denote the value of the t -th iterate of an optimization variable, the alternating minimization algorithm updates

$$\mu^{(t+1)} \leftarrow \arg \min_{\mu \in \mathbb{R}^{(r \times c)J}} \mathcal{F}_0(\mu, \Theta^{(t)}, \Gamma^{(t)}), \quad (9)$$

$$\Theta^{(t+1)} \leftarrow \arg \min_{\Theta \in \mathbb{R}^{(r \times c)J(J-1)/2}} \mathcal{F}_\rho(\mu^{(t+1)}, \Theta, \Gamma^{(t)}), \quad (10)$$

$$\Gamma_{j,m}^{(t+1)} \leftarrow \Gamma_{j,m}^{(t)} + \rho \left(\Theta_{j,m}^{(t+1)} - \mu_j^{(t+1)} + \mu_m^{(t+1)} \right) \text{ for } 1 \leq j < m \leq J,$$

until convergence. The ADMM algorithm modifies (9) by using \mathcal{F}_ρ rather than \mathcal{F}_0 . The advantage of using \mathcal{F}_0 is that we avoid solving an $rc \times rc$ linear system of equations at complexity $O(r^2c^2)$ when using the Kronecker structure. Using \mathcal{F}_0 also allows the updates for μ_1, \dots, μ_J to be computed in parallel with closed form solutions for each. Two conditions for the convergence of alternating minimization are that g is strongly convex (Tseng, 1991), which it is in our case, and that ρ is sufficiently close to zero. We provide a computable bound on the step size ρ to ensure convergence of our alternating minimization algorithm in the subsequent section.

The computational advantage of alternating minimization over ADMM was also recognized by Chi and Lange (2015) in the context of convex clustering. They found that the simplification of (9) relative to the ADMM version yielded a substantially more efficient algorithm.

Using the first order optimality condition for (9),

$$\mu_j^{(t+1)} = \bar{x}_j + \frac{1}{2\hat{\pi}_j} \Phi^{-1} \left(\sum_{\{m:m>j\}} \Gamma_{j,m}^{(t)} - \sum_{\{m:m<j\}} \Gamma_{m,j}^{(t)} \right) \Delta^{-1} \quad j = 1, \dots, J, \quad (11)$$

where $\hat{\pi}_j = n_j/n$ for $j = 1, \dots, J$.

The zero subgradient equation for (10) is

$$\rho \Theta_{j,m}^{(t+1)} + \Gamma_{j,m}^{(t)} - \rho \left(\mu_j^{(t+1)} - \mu_m^{(t+1)} \right) + \left\{ \lambda_1 w_{j,m} \circ h \left(\Theta_{j,m}^{(t+1)} \right) \right\} = 0, \quad (12)$$

where $h : \mathbb{R}^{r \times c} \rightarrow \mathbb{R}^{r \times c}$ and for all $(s, t) \in \{1, \dots, r\} \times \{1, \dots, c\}$,

$$[h(x)]_{s,t} = \begin{cases} \text{sign}(x_{s,t}) & : x_{s,t} \neq 0 \\ [-1, 1] & : x_{s,t} = 0 \end{cases}.$$

Tibshirani (1996), among others, have shown that (12) can be solved using the soft-thresholding

operator: $\text{soft}(x, \tau) = \max(|x| - \tau, 0)\text{sign}(x)$. The update for $\Theta_{j,m}$ is

$$\Theta_{j,m}^{(t+1)} = \text{soft} \left(\mu_j^{(t+1)} - \mu_m^{(t+1)} - \rho^{-1} \Gamma_{j,m}^{(t)}, \frac{\lambda_1}{\rho} w_{j,m} \right),$$

where soft is applied elementwise.

We use an accelerated variation of the algorithm presented in this section to solve (7). This is based on Goldstein et al. (2014) with simple restarting rules described by O'Donoghue and Candes (2015). Further details about our implementation are given in the subsequent section.

3.4 Summary

The block-wise coordinate descent algorithm for solving (3) is presented in Algorithm 1.

Algorithm 1. Given $\epsilon > 0$, $\Delta^{(0)} \in \mathbb{S}_c^+$, $\Phi^{(0)} \in \mathbb{S}_r^+$ such that $\|\Phi^{(0)}\|_1 = r$. Set $m = 0$:

Step 1: Compute $\mu^{(m+1)} = \arg \min_{\mu \in \mathbb{R}^{(r \times c)J}} g(\mu, \Phi^{(m)}, \Delta^{(m)}) + \lambda_1 \sum_{j < m} \|w_{j,m} \circ (\mu_j - \mu_m)\|_1$ using the algorithm from Section 3.3.

Step 2: Compute $\tilde{\Delta} = \text{GL} \{S_\delta(\mu^{(m+1)}, \Phi^{(m)}), \lambda_2\}$.

Step 3: Compute $\tilde{\Phi} = \text{GL} \{S_\phi(\mu^{(m+1)}, \tilde{\Delta}), \frac{\lambda_2}{c} \|\tilde{\Delta}\|_1\}$.

Step 4: Compute $\Delta^{(m+1)} = \frac{\|\tilde{\Phi}\|_1}{r} \tilde{\Delta}$, $\Phi^{(m+1)} = \frac{r}{\|\tilde{\Phi}\|_1} \tilde{\Phi}$

Step 5: If $f(\mu^{(m)}, \Phi^{(m)}, \Delta^{(m)}) - f(\mu^{(m+1)}, \Phi^{(m+1)}, \Delta^{(m+1)}) < \epsilon |f(\bar{x}, \Phi^{(0)}, \Delta^{(0)})|$, then stop. Otherwise, replace m by $m + 1$ and go to step 1.

In our implementation, we set $\epsilon = 10^{-6}$. To get initial values $\Phi^{(0)}$ and $\Delta^{(0)}$, we run the maximum likelihood algorithm (Dutilleul, 1999) until a mild convergence tolerance is reached, and use $\Phi^{(0)} = \text{diag}(\Phi^{\text{MLE}})$ and $\Delta^{(0)} = \text{diag}(\Delta^{\text{MLE}})$ where $(\Phi^{\text{MLE}}, \Delta^{\text{MLE}})$ are the final iterates.

Let $k_\phi^{(m)} = \varphi_{\min}(\Phi^{(m)})$ and $k_\delta^{(m)} = \varphi_{\min}(\Delta^{(m)})$, where φ_{\min} denotes the minimum eigenvalue. For the $(m + 1)$ th update of μ , if we select the step size parameter

$$\rho^{(m+1)} \in \left(0, \left\{ \min_j \{\hat{\pi}_j\} 4k_\phi^{(m)} k_\delta^{(m)} \right\} / J \right), \quad (13)$$

then the alternating minimization algorithm converges (Tseng, 1991; Chi and Lange, 2015). One can verify that (7) and (13) satisfy the conditions for convergence stated in section 6.2 of the supplemental material of Chi and Lange (2015) using an argument similar to theirs. The minimum eigenvalues of $\Phi^{(m)}$ and $\Delta^{(m)}$ are positive as long as initializers $\Phi^{(0)}$ and $\Delta^{(0)}$ are positive definite. When $k_\delta^{(m)}$ and $k_\phi^{(m)}$ are positive, g is strongly convex in μ , which is required for convergence.

In practice, we find it better to use ρ an order of magnitude smaller than the upper bound in (13), i.e., we use $\rho^{(m+1)} = (\min_j \{\hat{\pi}_j\} 4k_\phi^{(m)} k_\delta^{(m)}) / (10J)$ to ensure numerical stability. Although the step size $\rho^{(m+1)}$ may be small when $\Phi^{(m)}$ and $\Delta^{(m)}$ are dense, we find that when using accelerations and warm-starts, the small step size is not problematic.

We use an accelerated version of the alternating minimization algorithm proposed by Goldstein et al. (2014), which was also used by Chi and Lange (2015). O’Donoghue and Candes (2015) showed that acceleration restarts imposed after a fixed number of iterations can decrease the number of iterations required for convergence. In our implementation of the alternating minimization algorithm, we restart the accelerations after 200 iterations. We warm-start the $(m + 1)$ th update of μ by initializing the Lagrangian variables at their final iterates from the m th update.

At convergence of the alternating minimization algorithm, zeros in the final iterate of $\Theta_{j,m}$ do not correspond to exact entrywise equality in the final iterates for μ_j and μ_m . To enforce equality at the solution, we use simple thresholding.

3.5 Computational complexity

Solving (3) with $\lambda_1 = \lambda_2 = 0$, i.e. maximum likelihood estimation, also requires a blockwise coordinate descent algorithm (Dutilleul, 1999). The maximum-likelihood blockwise coordinate descent algorithm has computational complexity of order $O(nr^2c + nc^2r + r^3 + c^3)$. The first two terms come from computing the sample covariance matrices S_ϕ and S_δ , and the last two terms come from inverting S_ϕ and S_δ .

Our algorithm’s computational complexity is also $O(nr^2c + nc^2r + r^3 + c^3)$. We compute S_ϕ and S_δ and the graphical-lasso algorithm that we use is known to have worst case complexity $O(p^3)$ for a estimating a $p \times p$ precision matrix (Witten et al., 2011). In addition, for each μ update, we compute eigendecompositions of the iterates for Φ and Δ . The alternating minimization algorithm costs $O(r^2c + c^2r)$ when implemented in parallel.

The magnitude of tuning parameters effects the computing time of our algorithm. Generally, smaller values of λ_2 take longer.

4 Simulation study

4.1 Models

For 100 independent replications, we generated a realization of $n = n_{\text{train}} + n_{\text{validate}} + n_{\text{test}}$ independent copies of (X, Y) , where we set $n_{\text{train}} = n_{\text{validate}} = 75$, and $n_{\text{test}} = 1000$. The categorical response Y has support $\{1, 2, 3\}$ with probabilities $\pi_{*1} = \pi_{*2} = \pi_{*3} = 1/3$. Then

$$\text{vec}(X) \mid Y = j \sim N_{rc} \{ \text{vec}(\mu_{*j}), \Sigma_* \},$$

where μ_{*1} , μ_{*2} , and μ_{*3} are only different in one 4×4 submatrix, whose position is chosen randomly in each replication. We used multiple choices for the entries in this submatrix, which are displayed in Figure 1. All other mean matrix entries were set to zero. We consider four covariance models:

- Model 1. $\Sigma_* = \Delta_* \otimes \Phi_*$ where Φ_* has (a, b) th entry $0.7^{|a-b|}$ and Δ_* has (c, d) th entry $0.7 \times 1(c \neq d) + 1(c = d)$.
- Model 2. $\Sigma_* = \Delta_* \otimes \Phi_*$ where Φ_* has (a, b) th entry $0.7^{|a-b|}$ and Δ_* is block-diagonal where

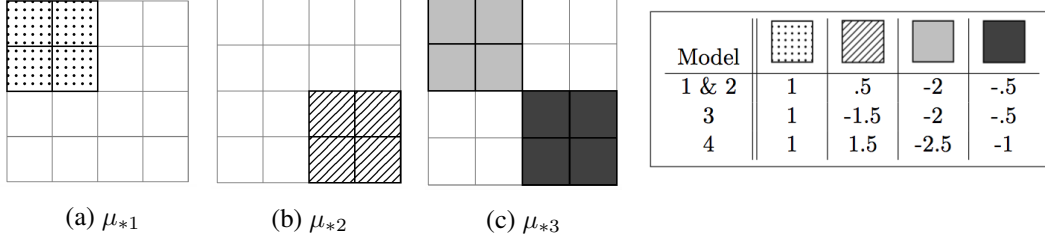


Figure 1: The 4×4 submatrix where μ_{*1} , μ_{*2} , and μ_{*3} differ. White corresponds to zero and the legend gives the values corresponding to the highlighted cells for each model.

Δ_* can be expressed elementwise:

$$\Delta_{c,d} = \begin{cases} 1 & \text{if } c = d \\ 0.7 & \text{if } \mu_{*j,a,c} \neq \mu_{*m,a,d} \text{ for any } a \in \{1, \dots, r\} \text{ and } 1 \leq j < m \leq J \\ 0 & \text{otherwise} \end{cases}$$

- Model 3. Σ_* corresponds to the covariance model

$$\text{Cov}(X_{a,b}, X_{c,d} \mid Y = j) = \{0.5I(b \neq d) + I(b = d)\} \frac{(\rho_b \rho_d)^{|a-c|}}{1 - \rho_b \rho_d},$$

where ρ_1, \dots, ρ_c are c equally spaced values between 0.5 and 0.9. The matrix Σ_* is positive definite when $r = c$ with $c = \{8, 16, 32, 64\}$, and when $r = 32$ with $c = \{8, 16, 32, 64\}$.

- Model 4. Σ_* corresponds to the covariance model

$$\text{Cov}(X_{a,b}, X_{c,d} \mid Y = j) = \begin{cases} 1 & \text{if } (a, b) = (c, d) \\ 0.5 & \text{if } \mu_{*j,a,b} \neq \mu_{*m,c,d} \text{ for any } 1 \leq j < m \leq J \\ 0 & \text{otherwise} \end{cases}$$

In Model 3, if $\rho_k = \rho$ for all $k \in \{1, \dots, c\}$, then Σ_* has the decomposition (2) corresponding to Φ_* with an AR(1) structure and Δ_* with a compound symmetric structure (Mitchell et al., 2006). However, when $\rho_k \neq \rho$, Σ_* does not have decomposition (2): the covariance between any two entries in the same row depends on the column and vice versa. Model 4 is the rc -variate normal model similar to the first model used in the simulations from Xu et al. (2015).

4.2 Methods

We consider the following model-based methods for fitting the linear discriminant analysis model:

- Bayes. The Bayes rule, i.e., Σ_* , μ_* , and π_{*j} known for $j = 1, \dots, J$;
- MN. The maximum likelihood estimator of (1) under (2), i.e., the matrix-normal maximum likelihood estimator;
- Guo. The sparse naïve Bayes type-estimator proposed by Guo (2010) defined in Section 2.2 with tuning parameter chosen to minimize misclassification rate on the validation set;

- **vec-SURE.** The multiclass SURE independence screening method proposed by Pan et al. (2016) with model sizes chosen to minimize misclassification rate on the validation set;
- **MN-SURE.** The matrix-normal extension of the SURE independence screening estimator proposed by Pan et al. (2016) with model sizes chosen to minimize misclassification error on the validation set.
- **PMN(μ).** The estimator defined by (3) with $\mu = \mu_*$ fixed and λ_2 chosen to minimize misclassification rate on the validation set;
- **PMN(Σ) / Xu(Σ).** The estimator defined by (3) with $\Phi = \Phi_*$ and $\Delta = \Delta_*$ fixed when $\Sigma_* = \Delta_* \otimes \Phi_*$; the estimator defined by (4) with $\hat{\Sigma} = \Sigma_*$ fixed when $\Sigma_* \neq \Delta_* \otimes \Phi_*$; and λ_1 chosen to minimize misclassification rate on the validation set;
- **PMN.** The estimator defined by (3) with tuning parameters chosen by minimizing misclassification rate on the validation set.

The methods PMN(μ) and PMN(Σ) / Xu(Σ) both use some oracle information and were included to study how estimating μ_* , Δ_* , and Φ_* simultaneously affect classification accuracy. We refer to these method as part-oracle matrix-LDA methods. We refer to Guo and vec-SURE as vector-LDA methods; MN and MN-SURE as non-oracle matrix-LDA methods. MN-SURE is a matrix-normal generalization of the screening method proposed by Pan et al. (2016).

Following Guo (2010), we use a validation set to select tuning parameters. The candidate set for tuning parameters was $\{2^x : x = -12, -11.5, \dots, 11.5, 12\}$. Candidate model sizes for vec-SURE and MN-SURE were $\{0, 1, \dots, 25\}$, where model size refers to the number of pairwise nonzero mean differences based on thresholding.

4.3 Performance measures

To compare classification accuracy, we record the misclassification rate on the test set for each replication. We also measure identification of mean differences that are zero through both true positive rate (TPR) and true negative rate (TNR). Let $D(\mu_*) = [\text{vec}(\mu_{*1} - \mu_{*2}), \dots, \text{vec}(\mu_{*(J-1)} - \mu_{*J})]$, and $D(\hat{\mu}) = [\text{vec}(\hat{\mu}_1 - \hat{\mu}_2), \dots, \text{vec}(\hat{\mu}_{(J-1)} - \hat{\mu}_J)]$. We define TPR as

$$\text{TPR}(\hat{\mu}, \mu_*) = \frac{\# \left\{ (z, w) : [D(\hat{\mu})]_{z,w} \neq 0 \cap [D(\mu_*)]_{z,w} \neq 0 \right\}}{\# \left\{ (z, w) : [D(\mu_*)]_{z,w} \neq 0 \right\}},$$

where $\#$ denotes cardinality. We similarly define TNR as

$$\text{TNR}(\hat{\mu}, \mu_*) = \frac{\# \left\{ (z, w) : [D(\hat{\mu})]_{z,w} = 0 \cap [D(\mu_*)]_{z,w} = 0 \right\}}{\# \left\{ (z, w) : [D(\mu_*)]_{z,w} = 0 \right\}}.$$

TNR and TPR together address mean difference estimation which we use as a measure of variable selection for comparison to the estimator of Guo (2010) and Pan et al. (2016).

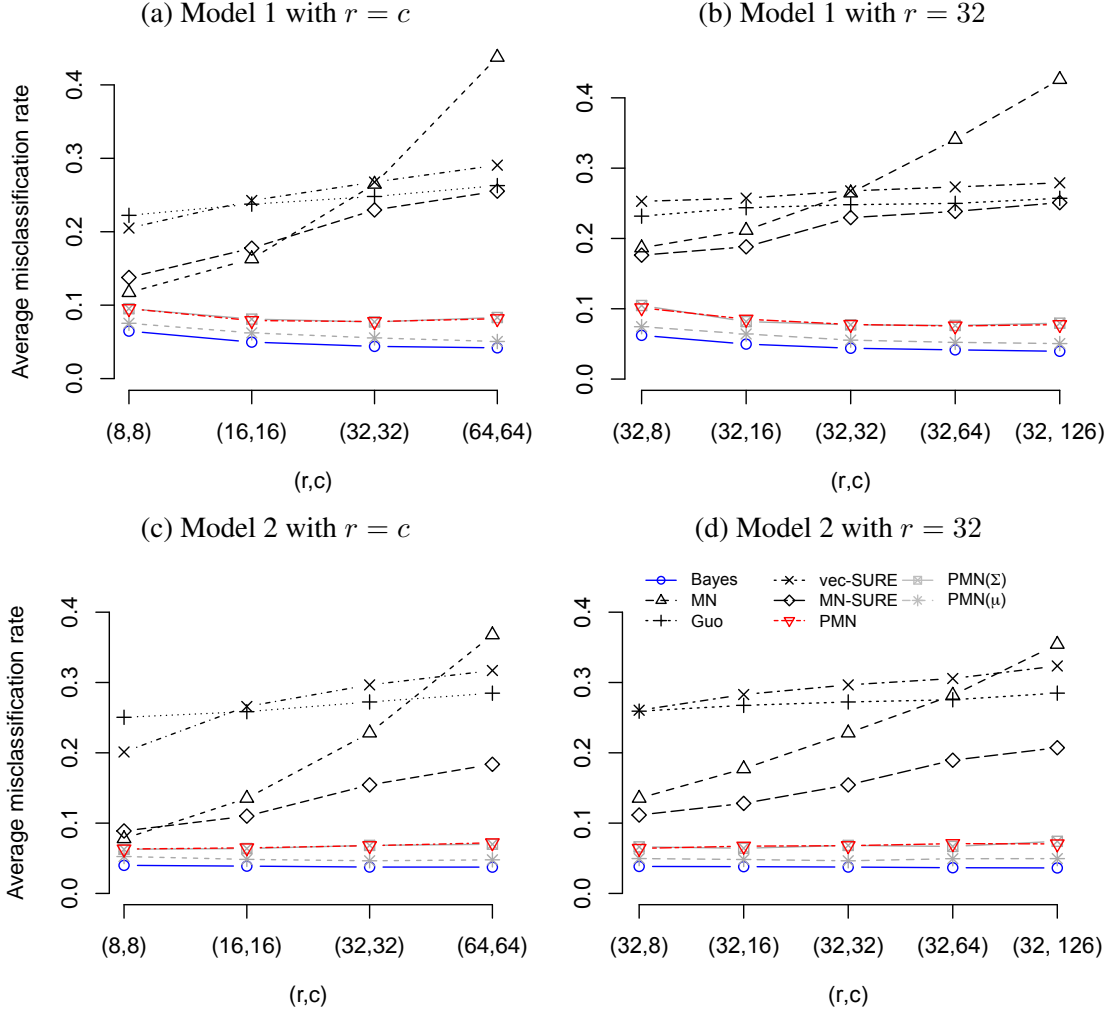


Figure 2: Misclassification rates averaged over 100 replications; (a) and (b) are for Model 1 and (c) and (d) for Model 2.

4.4 Results

We display average misclassification rates for Models 1 and 2 in Figure 2. For Model 1, the matrix-normal maximum likelihood estimator tended to outperform the vector-LDA methods when r and c were small, but its average classification rate got worse as the dimensionality increases. The estimator proposed by Guo (2010) performs poorly when r and c are small, but got worse more slowly than the other vector and non-oracle matrix-LDA methods. The misclassification rate of the Bayes rules suggests that as the dimensionality increases in Model 1, the optimal misclassification rate can be improved. Our method PMN had improved classification accuracy as both r and c increased and performed similarly to $\text{PMN}(\Sigma)$, which uses some oracle information.

TPR and TNR results are displayed in Table 1. For Model 1, PMN tended to have the second highest TNR behind vec-SURE, but tends to have higher TPR than all competing methods except $\text{PMN}(\Sigma)$, which uses some oracle information.

Results were similar for Model 2. The matrix-normal variation of the SURE screening esti-

Table 1: TNR/TPR percentages averaged over the 100 replications for Model 1-4.

Method	Model 1 (r, c)							
	(8,8)	(16,16)	(32,32)	(64,64)	(32,8)	(32,16)	(32,64)	(32,126)
Guo	85.7/79.4	95.8/68.8	98.6/65.6	99.4/59.8	96.8/70.9	97.6/68.2	99.2/65.5	99.5/59.4
vec-SURE	88.5/71.9	97.9/52.4	99.7/40	99.8/35.4	98.4/49.9	99.2/48.1	99.8/38.9	99.9/35.2
MN-SURE	35.2/90.9	80.2/66.9	97.3/46.5	99.2/37.9	87.1/64.8	90/61.5	98.4/43.6	99.4/38.9
PMN(Σ)	85.9/88.2	94/84.1	98.2/78.4	99/81.2	94.1/82.6	96.7/83.5	98.7/80.6	99.2/74.9
PMN	95.1/79.9	95.8/77.5	99/74	99.5/69.9	98.2/74.2	98.6/74.6	99.3/73.9	99.6/71.6
	Model 2 (r, c)							
	(8,8)	(16,16)	(32,32)	(64,64)	(32,8)	(32,16)	(32,64)	(32,126)
Guo	81.4/81.6	94.7/73.2	97.5/65.8	98.6/60.8	94.6/74.1	96.6/68.8	99/62.5	99/63.1
vec-SURE	87.1/71.2	98.1/51.1	99.7/36.2	99.9/29.5	98.1/51.4	99.2/47	99.8/35.5	99.9/32.2
MN-SURE	46.1/89	74.2/72.9	91.8/54	97.6/42.9	83.9/68.9	86.4/65.5	96.9/43.9	98.2/39.4
PMN(Σ)	90.9/87.5	94.4/86.9	98.8/80.6	99.6/84.2	95/86.9	97.4/85.5	99.2/82.5	99.4/79.9
PMN	96.5/79.1	96.9/77.5	99.1/73	99.8/70.5	98.7/77.6	99.1/74.2	99.5/68.9	99.7/70.9
	Model 3 (r, c)							
	(8,8)	(16,16)	(32,32)	(64,64)	(32,8)	(32,16)	(32,64)	(32,126)
Guo	86.8/84.2	95.5/81.2	97.9/75.1	99.4/62.6	95.3/80.4	96/79.8	98.7/70.2	—/—
vec-SURE	84.8/82.9	97.2/65.5	99.4/44.6	99.8/31.2	97.8/55.5	98.8/52.6	99.7/37	—/—
MN-SURE	38.6/94.2	80.6/81.5	96.4/53.2	98.8/35.5	85.5/72.8	90.8/67.4	97.7/43	—/—
Xu(Σ)	81.4/92.1	93.8/90	96.3/86.9	98.5/83.9	91.4/87.1	95.7/87.9	97.9/87.1	—/—
PMN	86.4/96.1	93.8/96	98.3/93	99.3/87.2	93.8/93.5	95.7/94.1	98.8/90.1	—/—
	Model 4 (r, c)							
	(8,8)	(16,16)	(32,32)	(64,64)	(32,8)	(32,16)	(32,64)	(32,126)
Guo	82.2/98.5	93.9/98.5	96.9/97.1	98.9/96.1	90.4/98.2	95/97.9	97.8/93.4	98.9/96
vec-SURE	97.7/83.1	99.4/79.9	99.8/78.4	99.9/70.6	99.2/80.1	99.7/79.1	99.9/74.6	99.9/74.9
MN-SURE	73.1/97.1	89.8/96.2	96/91	98.8/82.4	89.9/94.8	95.7/91.6	98.4/84.5	99/84.4
Xu(Σ)	87.7/99.2	93.3/98.4	97.8/96.9	99.5/97.1	92/97.8	96.3/96.2	98.8/95.8	99.4/96.1
PMN	92.6/96.6	95.8/97.5	97.3/94.9	99.5/92.2	94.3/96.4	97.6/95.5	99.2/91.5	99.4/93.8

mator of Pan et al. (2016) tended to perform best among the vector and non-oracle matrix-LDA methods. The estimator of Guo (2010) got worse the slowest amongst the vector-LDA methods. PMN performed as well as PMN(Σ), both of which performed more closely to PMN(μ) and the Bayes rule than for Model 1.

The misclassification rates for Models 3 and 4 are displayed in Figure 3. In Model 3, although Σ_* does not have the Kronecker decomposition in (2), PMN outperformed all non part-oracle estimators. In terms of TPR and TNR results presented in Table 3, PMN performed similarly to Xu(Σ), both of which had higher TPR than competitors and TNR similar to vec-SURE. This suggests that even when (2) does not hold, our method can perform well in classification.

In Model 4, PMN performed similarly to the vector-LDA methods. MN-SURE was the best non-oracle method, which suggests that (2) may be a reasonable alternative to naïve Bayes under high dimensionality. Like in Model 3, PMN performed similarly to Xu(Σ) in terms of TPR and TNR.

5 EEG data example

We analyzed the EEG data (<https://kdd.ics.uci.edu/databases/eeg/eeg.html>) also studied by Li et al. (2010) and Zhou and Li (2014). In the original study, 122 subjects, 77 of whom were alcoholics and 45 of whom were control, were exposed to stimuli while voltage was measured from $c = 64$

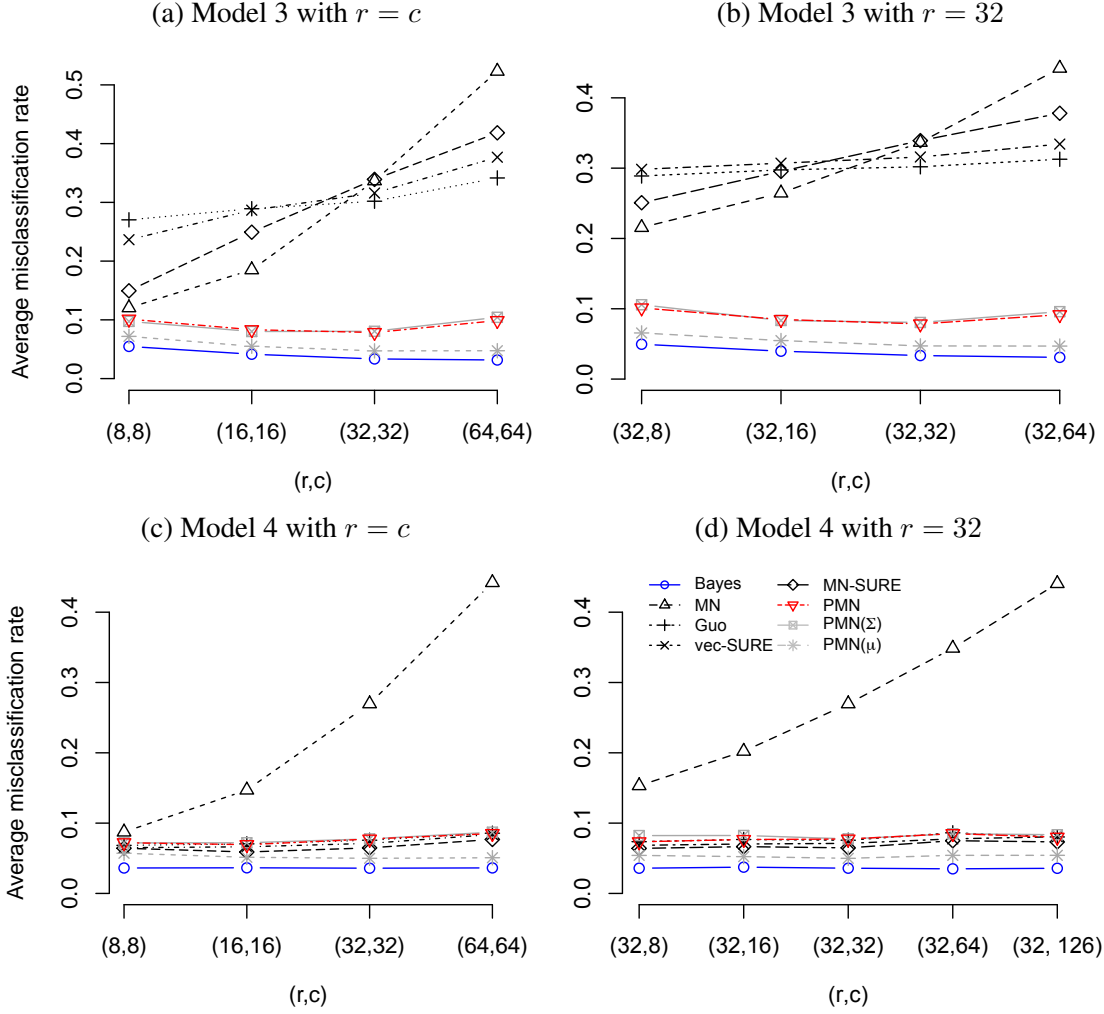


Figure 3: Misclassification rates averaged over 100 replications; (a) and (b) are for Model 3 and (c) and (d) for Model 4.

channels on a subject's scalp at $r = 256$ time points. Each subject underwent 120 trials. Each trial had one of three possible stimuli: single stimulus, two matched stimuli, or two unmatched stimuli. As in Li et al. (2010) and Zhou and Li (2014), we only analyze the single stimulus condition. Because each subject underwent multiple trials under the single stimulus condition, we use the within subject average over all single stimulus trials as the predictor and we use whether they were alcoholic or control as the response.

It is common to assume that (2) holds in the analysis of EEG data. For example, Zhou (2014) assumed that (2) holds when analyzing a single subject from this same dataset. It may also be reasonable to assume that only a subset of channels and time point combinations are important for discriminating between alcoholic and control response categories. Thus, the primary goal of our analysis is to identify a subset of channels and time point combinations that help explain how the alcoholics and controls react to the stimulus differently.

To demonstrate our method's classification accuracy, we used the leave-one-out cross validation approach from Li et al. (2010) and Zhou and Li (2014). For $k = 1, \dots, 122$, we left out the k th

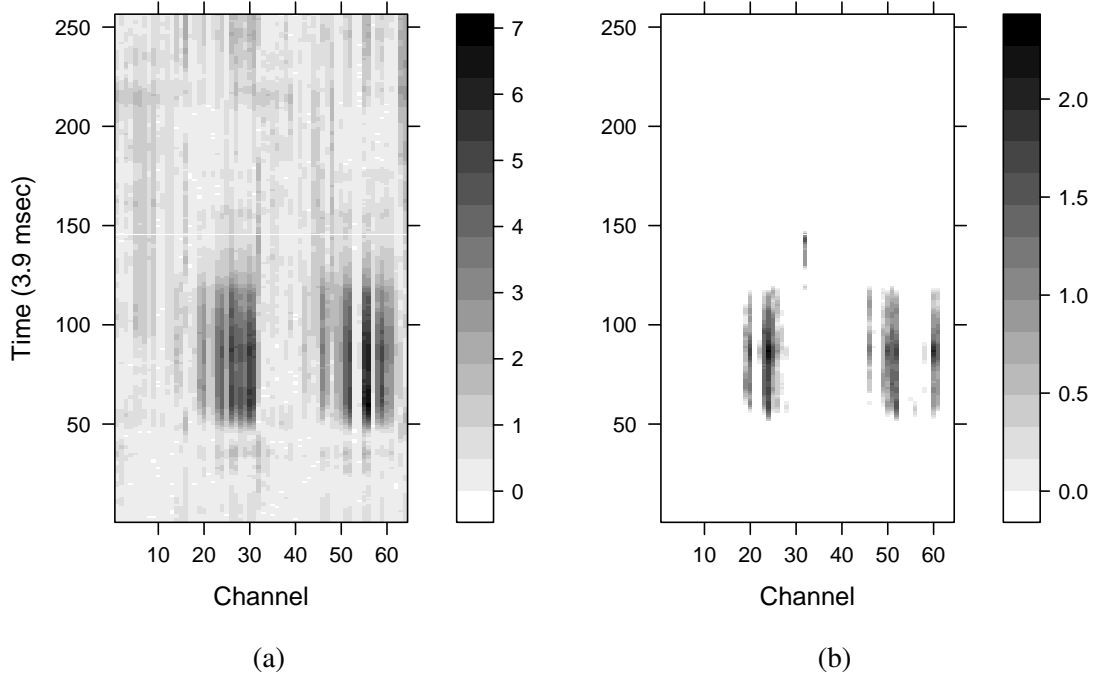


Figure 4: (a) The absolute value of the sample mean differences between the alcoholic and control response categories. (b) The absolute value of the estimated mean differences from (3) based on the tuning parameter pair $(\lambda_1, \lambda_2) = (0.15, 5.66)$, which had leave-one-out cross-validation classification accuracy of 98 out of 122.

observation and used the remaining 121 observations as training data. For each k , we selected tuning parameters for use in (3) by minimizing 5-fold cross validation misclassification error on the training dataset. Our method correctly classified 97 of 122 observations. Li et al. (2010) and Zhou and Li (2014) reported correctly classifying 97 and 94 of 122, respectively. Li et al. (2010) used quadratic discriminant analysis after dimension-folding of the predictors, and Zhou and Li (2014) used logistic regression with spectral regularization of the coefficient matrix.

To demonstrate the interpretability our fitted model, we separately fit (3) using the complete dataset. We used a tuning parameter pair $(\lambda_1, \lambda_2) = (0.15, 5.66)$, which had leave-one-out classification accuracy of 98 out of 122. The estimated mean difference, displayed as a heatmap in Figure 4(b), had 15466 of 16384 entries equal to zero.

Our fitted model can be used to easily identify which channels and time points have nonzero mean differences. We estimated only 22 of the 64 channels to have at least one time point where the mean differences were nonzero, only 16 of which had at least 20 nonzero time points. Inspecting the estimated mean differences displayed in Figure 4, it seems that the majority of activity that distinguishes between the alcoholic and control subjects takes place between the 52nd and 115th time points. We used the R package *eegkit* (Helwig, 2015) to display which channels had nonzero mean differences in Figure 4a. Our method does not explicitly use the spatial structure of channels in estimation, yet it recovered a set of important channels which have a natural arrangement in space.

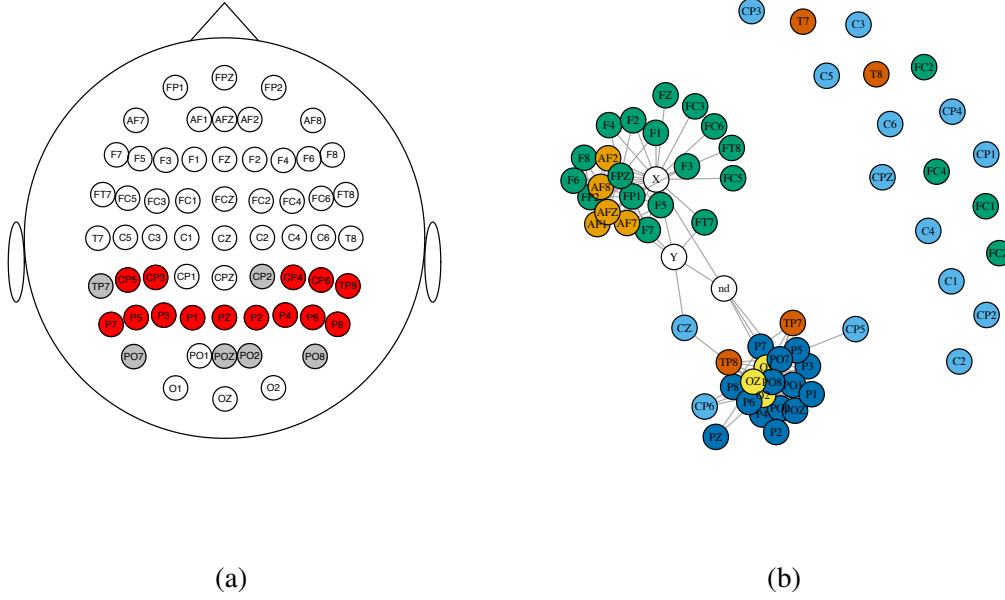


Figure 5: (a) An EEG cap based on the fitted model using $(\lambda_1, \lambda_2) = (0.15, 5.66)$. Red channels had at least 20 time points estimated to have nonzero mean differences; grey channels had less than 20 but greater than zero, whereas white channels had no nonzero mean differences. (b) The Gaussian precision graphical model associated with $\hat{\Delta}$. Colors correspond to different regions of the EEG channels; white channels are those that do not appear on the EEG cap image.

Both Φ_* and Δ_* were estimated to be relatively sparse: $\hat{\Phi}$ was a diagonal matrix, while $\hat{\Delta}$ had 3676 of 4032 off-diagonals equal to zero. Our estimate $\hat{\Delta}$ can be interpreted in terms of a Gaussian precision graphical model corresponding to the conditional dependence structure of the channels. We display the graphical model corresponding to $\hat{\Delta}$ in Figure 5b. The graph structure corresponds to the spatial arrangement of channels displayed in Figure 5a – a result also observed by Zhou (2014).

6 Extension to quadratic discriminant analysis

Our method naturally extends to the quadratic discriminant analysis model, where one assumes

$$\text{vec}(X) \mid Y = j \sim N_{rc} \{ \text{vec}(\mu_{*j}), \Sigma_{*j} \}, \quad j = 1, \dots, J$$

where $\Sigma_{*j} \in \mathbb{S}_{rc}^+$ is the covariance matrix for the j th response category. To generalize (2), one can assume either (i) $\Sigma_{*j}^{-1} = \Delta_{*j} \otimes \Phi_{*j}$, (ii) $\Sigma_{*j}^{-1} = \Delta_{*j} \otimes \Phi_*$, or (iii) $\Sigma_{*j}^{-1} = \Delta_* \otimes \Phi_{*j}$. Our algorithms can be modified to accommodate these cases.

Acknowledgments

This research was supported in part by the Doctoral Dissertation Fellowship from the University of Minnesota and the National Science Foundation grant DMS-1452068.

References

- Allen, G. I. and Tibshirani, R. (2010). Transposable regularized covariance models with an application to missing data imputation. *The Annals of Applied Statistics*, 4(2):764.
- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine Learning*, 3(1):1–122.
- Chi, E. C. and Lange, K. (2015). Splitting methods for convex clustering. *Journal of Computational and Graphical Statistics*, 24(4):994–1013.
- Dutilleul, P. (1999). The mle algorithm for the matrix normal distribution. *Journal of Statistical Computation and Simulation*, 64(2):105–123.
- Friedman, J., Hastie, T., and Tibshirani, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441.
- Goldstein, T., O’Donoghue, B., Setzer, S., and Baraniuk, R. (2014). Fast alternating direction optimization methods. *SIAM Journal on Imaging Sciences*, 7(3):1588–1623.
- Guo, J. (2010). Simultaneous variable selection and class fusion for high-dimensional linear discriminant analysis. *Biostatistics*, 11(4):599–608.
- Gupta, A. K. and Nagar, D. K. (2000). *Matrix variate distributions*. Chapman and Hall/CRC Press.
- Helwig, N. E. (2015). *eegkit: Toolkit for Electroencephalography Data*. R package version 1.0-2.
- Hung, H. and Wang, C.-C. (2013). Matrix variate logistic regression model with application to eeg data. *Biostatistics*, 14(1):189–202.
- Hunter, D. R. and Li, R. (2005). Variable selection using mm algorithms. *The Annals of Statistics*, 33(4):1617.
- Leng, C. and Tang, C. Y. (2012). Sparse matrix graphical models. *Journal of the American Statistical Association*, 107(499):1187–1200.
- Li, B., Kim, M. K., and Altman, N. (2010). On dimension folding of matrix-or array-valued statistical objects. *The Annals of Statistics*, 38(2):1094–1121.
- Li, M. and Yuan, B. (2005). 2d-lda: A statistical linear discriminant analysis for image matrix. *Pattern Recognition Letters*, 26(5):527–532.
- Mitchell, M. W., Genton, M. G., and Gumpertz, M. L. (2006). A likelihood ratio test for separability of covariances. *Journal of Multivariate Analysis*, 97(5):1025–1043.

- O'Donoghue, B. and Candes, E. (2015). Adaptive restart for accelerated gradient schemes. *Foundations of Computational Mathematics*, 15(3):715–732.
- Pan, R., Wang, H., and Li, R. (2016). Ultrahigh-dimensional multiclass linear discriminant analysis by pairwise sure independence screening. *Journal of the American Statistical Association*, 111(513):169–179.
- Roś, B., Bijma, F., de Munck, J. C., and de Gunst, M. C. (2016). Existence and uniqueness of the maximum likelihood estimator for models with a kronecker product covariance structure. *Journal of Multivariate Analysis*, 143:345–361.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B*, 58:267–288.
- Tseng, P. (1991). Applications of a splitting algorithm to decomposition in convex programming and variational inequalities. *SIAM Journal on Control and Optimization*, 29(1):119–138.
- Tsiligkaridis, T., Hero, A. O., and Zhou, S. (2012). Kronecker graphical lasso. In *2012 IEEE Statistical Signal Processing Workshop (SSP)*, pages 884–887. IEEE.
- Witten, D. M., Friedman, J. H., and Simon, N. (2011). New insights and faster computations for the graphical lasso. *Journal of Computational and Graphical Statistics*, 20(4):892–900.
- Xu, P., Zhu, J., Zhu, L., and Li, Y. (2015). Covariance-enhanced discriminant analysis. *Biometrika*, 102(1):33–45.
- Zhang, Y. and Schneider, J. G. (2010). Learning multiple tasks with a sparse matrix-normal penalty. In *Advances in Neural Information Processing Systems*, pages 2550–2558.
- Zhong, W. and Suslick, K. S. (2015). Matrix discriminant analysis with application to colorimetric sensor array data. *Technometrics*, 57(4):524–534.
- Zhou, H. and Li, L. (2014). Regularized matrix regression. *Journal of the Royal Statistical Society: Series B*, 76(2):463–483.
- Zhou, S. (2014). Gemini: Graph estimation with matrix variate normal instances. *The Annals of Statistics*, 42(2):532–562.