

Statistica Sinica Preprint No: SS-2020-0496

Title	Parsimonious Tensor Discriminant Analysis
Manuscript ID	SS-2020-0496
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.202020.0496
Complete List of Authors	Ning Wang, Wenjing Wang and Xin Zhang
Corresponding Authors	Xin Zhang
E-mails	henry@stat.fsu.edu

Parsimonious Tensor Discriminant Analysis

Ning Wang, Wenjing Wang, and Xin Zhang

Florida State University, Tallahassee, FL, 32306

Abstract: Discriminant analyses of multidimensional array data (i.e., tensors) are of substantial interest in numerous statistics and engineering research problems, such as signal processing, imaging, genetics, and brain–computer interfaces. In this study, we consider a multi-class discriminant analysis with a tensor-variate predictor and a categorical response. To overcome the high dimensionality and to exploit the tensor correlation structure, we propose the discriminant analysis with tensor envelope (DATE) model for simultaneous dimension reduction and classification. We extend the notion of tensor envelopes from regression to discriminant analysis and develop two complementary estimation procedures: DATE-L is a likelihood-based estimator that is shown to be asymptotically efficient when the sample size goes to infinity and the tensor dimension is fixed; DATE-D is a novel decomposition-based estimator suitable for high-dimensional problems. Interestingly, we show that DATE-D is still root-n consistent, even when the tensor dimensions on each model grow arbitrarily fast, but at a similar rate. We demonstrate the robustness and efficiency of our estimators using extensive simulations and real-data examples.

Key words and phrases: Dimension reduction; Linear discriminant analysis; Tensor.

The authors contributed equally to this work and are listed in alphabetical order. The authors would like to thank the associate editor and three reviewers for their helpful comments. This research was supported in part by grants CCF-1908969, DMS-2053697, and DMS-2113590 from the U.S. National Science Foundation.

1. Introduction

Statistical analyses of tensor data are common in areas such as high-throughput genetics (Hore et al. 2016), signal processing (Cichocki et al. 2015), neuroimaging (Zhou et al. 2013), and point cloud data (Yan et al. 2019), among others. Our notion of a tensor analysis differs from that in mathematics and physics, although some operators and techniques are the same. We use multilinear algebra (e.g., Hitchcock 1927, Tucker 1966) to provide a concise statistical modeling framework and estimation procedures.

Two tasks popular in the statistical literature on tensor data analysis are tensor decompositions and tensor regression problems. In the first category, studies are mostly unsupervised, with tensor decompositions used to reduce the dimensionality of the tensor and to extract a meaningful representation along each mode of the tensor. For example, Kolda & Bader (2009) give an overview of tensor decompositions and related applications, Wang et al. (2018) provide a recent review on tensor sparse recovery, Chi & Kolda (2012) developed algorithms for sparse count data, Zhang & Xia (2018) study the theoretical limits of tensor singular value decomposition, and Zhang (2019) examine tensor completion. In the second category, the goal is to study the relationship between a tensor variable and other variables (scalar, vector, or even tensor). Such problems are formulated as tensor regression problems. In particular, a tensor variable can appear in regression models as the predictor (e.g., Zhou et al. 2013, Wang et al. 2017, Li et al. 2018), the response (e.g., Hoff 2015, Li & Zhang 2017, Sun &

1. INTRODUCTION³

Li 2017), or both (i.e., tensor on tensor regression, Lock 2018, Gahrooei et al. 2020, Raskutti et al. 2019).

Here, we study the problem of tensor discriminant analysis. Unlike the abundant literature on tensor decompositions and regression problems, there are far fewer statistical approaches and theoretical studies for tensor classification and discriminant analysis. Some earlier works (Ye et al. 2004, Li & Schonfeld 2014, Zhong & Suslick 2015, Yan et al. 2006) have shown promising performance for matrix and tensor discriminant analysis, based on the principle of maximizing the ratio of between-class variation to within-class variation. These methods are thus extensions of Fisher's discriminant analysis (Fisher 1936) to matrix/tensor data. Because of the high dimensionality, such linear/multilinear classifiers are arguably more reliable than quadratic or nonlinear discriminant analysis. Another important, but different research direction is using margin-based classification methods for tensor data (Lyu et al. 2017, Li & Maiti 2019). More recently, likelihood-based matrix/tensor discriminant analysis models and methods (Molstad & Rothman 2019, Pan et al. 2019) were shown to be more effective than moment-based and margin-based methods.

We propose the discriminant analysis with tensor envelope (DATE) model, which incorporates the recently proposed tensor envelope (Li & Zhang 2017, Zhang & Li 2017) to reduce the model complexity and improve estimation efficiency. A Tensor envelope is a multilinear extension of the envelope methodology in multivariate statistics (Cook et al. 2010). We provide a brief review of envelopes and tensor envelopes in Section 2.2. The core idea of an envelope

1. INTRODUCTION⁴

is to identify and eliminate unimportant and immaterial variation in the data in order to improve the estimation and prediction. This is achieved by projecting the data onto a latent subspace, known as an “envelope.” The existing envelope methods for tensor data were developed in regression problems, with very few focusing on tensor discriminant analysis. We address this gap in the literature by extending the envelope discriminant analysis (Zhang & Mai 2019) from vector data to tensor data. Similarly to existing envelope methods, we derive a likelihood-based estimator (DATE-L) that is asymptotically efficient. To accommodate high-dimensional applications, we propose a novel decomposition-based estimator (DATE-D) that is a complementary alternative to the more expensive manifold optimization in likelihood-based envelope estimation. We obtain a convergence rate for DATE-D that is sufficiently strong for most tensor data applications. Therefore, DATE-D provides a computationally feasible and theoretically justified approach in high dimensions when DATE-L fails. It also fills the gap in the literature on high-dimensional theoretical analysis of envelope methods, especially because additional structural assumptions, such as sparsity, are not required.

Extending the tensor envelope concept from regression to the present setting is not trivial, requiring new parameterization and maximum likelihood estimator derivations. We also adapt the fast and stable one-step estimation (Li & Zhang 2017) and 1D algorithm (Cook & Zhang 2016) to control the computational complexity of our DATE-L estimation procedure. More importantly, existing envelope methods (including DATE-L) often require iterative Grassmann man-

2. BACKGROUND

ifold optimization. We provide a novel decomposition-based estimation that is computationally tractable and theoretically justified for high-dimensional tensors, and can be straightforwardly modified to fit tensor envelope regression models in high dimensions. While existing tensor envelopes are studied under fixed tensor dimensions, we establish new consistency results for both fixed and diverging tensor dimensions.

The rest of the paper is structured as follows. In Section 2, we introduce some tensor notation and briefly review envelopes in both vector and tensor regression. In Section 3, we propose the DATE model and derive the two estimation procedures, DATE-L and DATE-D. Section 4 studies the asymptotic properties of the DATE-L estimator and the convergence rate of the DATE-D estimator. Simulations and real-data examples are presented in Section 5. Section 6 concludes the paper. Additional numerical results, implementation details, and proofs are provided in the online Supplementary Material.

2. Background

2.1 Notation

We call a multidimensional array $\mathbf{A} \in \mathbb{R}^{p_1 \times \cdots \times p_M}$ an M -way tensor or M th-order tensor (e.g., $M = 1$ for vectors and $M = 2$ for matrices). Some key operators on the M th-order tensor \mathbf{A} are defined as follows. **Vectorization:** The vectorization of \mathbf{A} is denoted by $\text{vec}(\mathbf{A}) \in \mathbb{R}^{\prod_{m=1}^M p_m}$, where the (i_1, \dots, i_M) th scalar in \mathbf{A} is mapped to the j th entry of $\text{vec}(\mathbf{A})$, for $j = 1 + \sum_{m=1}^M \{(i_m - 1) \prod_{k=1}^{m-1} p_k\}$.

2. BACKGROUND6

Matricization: The *mode- n matricization*, reshapes the tensor \mathbf{A} into a matrix denoted by $\mathbf{A}_{(n)} \in \mathbb{R}^{p_n \times \prod_{m \neq n} p_m}$, such that the (i_1, \dots, i_M) -th element in \mathbf{A} becomes the (i_n, j) th element of the matrix $\mathbf{A}_{(n)}$, where $j = 1 + \sum_{k \neq n} \{(i_k - 1) \prod_{l < k, l \neq n} p_l\}$. **(Collapsed) Vector product :** The *mode- n vector product* of \mathbf{A} and a vector $\mathbf{c} \in \mathbb{R}^{p_n}$ is represented by $\mathbf{A} \bar{\times}_n \mathbf{c} \in \mathbb{R}^{p_1 \times \dots \times p_{n-1} \times p_{n+1} \times \dots \times p_M}$, and results in a collapsed $(M - 1)$ th-order tensor. This product is the result of the inner products between every *mode- n fiber* in \mathbf{A} with vector \mathbf{c} . The *mode- n fibers* of \mathbf{A} are the vectors obtained by fixing all indices except the n th index. **Matrix product:** The *mode- n product* of a tensor \mathbf{A} and a matrix $\mathbf{G} \in \mathbb{R}^{s \times p_n}$, denoted as $\mathbf{A} \times_n \mathbf{G}$, is an M th-order tensor with dimension $p_1 \times \dots \times p_{n-1} \times s \times p_{n+1} \times \dots \times p_M$. Similarly to the vector product, the product is the result of a multiplication between each *mode- n fiber* of \mathbf{A} and the matrix \mathbf{G} . **Tucker product:** The *Tucker product* of the core tensor \mathbf{A} and a series of factor matrices $\mathbf{C}_1, \dots, \mathbf{C}_M$, where $\mathbf{C}_k \in \mathbb{R}^{q_k \times p_k}$, for $k = 1, \dots, m$, is defined as $\mathbf{A} \times_1 \mathbf{C}_1 \times_2 \dots \times_M \mathbf{C}_M \equiv \llbracket \mathbf{A}; \mathbf{C}_1, \dots, \mathbf{C}_M \rrbracket \in \mathbb{R}^{q_1 \times \dots \times q_M}$. If $q_k \geq p_k$ and each \mathbf{C}_k satisfies $\mathbf{C}_k^T \mathbf{C}_k = \mathbf{I}_{q_k}$, then $\mathbf{B} = \llbracket \mathbf{A}; \mathbf{C}_1, \dots, \mathbf{C}_M \rrbracket$ is called a Tucker decomposition of \mathbf{B} . **The tensor normal distribution** with mean $\boldsymbol{\mu} \in \mathbb{R}^{p_1 \times \dots \times p_M}$ and separable covariance matrices $\boldsymbol{\Sigma}_m > 0$, where $\boldsymbol{\Sigma}_m \in \mathbb{R}^{p_m \times p_m}$, for $m = 1, \dots, M$, is denoted by $TN(\boldsymbol{\mu}, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_M)$. We have $\mathbf{X} \sim TN(\boldsymbol{\mu}, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_M)$ if $\mathbf{X} = \boldsymbol{\mu} + \llbracket \mathbf{Z}; \boldsymbol{\Sigma}_1^{1/2}, \dots, \boldsymbol{\Sigma}_M^{1/2} \rrbracket$, where \mathbf{Z} is the standard tensor normal random variable with elements that are all independently $N(0, 1)$ random variables. The vectorization and matricization operations on a tensor normal random variable preserve the normality. Specifically, $\text{vec}(\mathbf{X})$

2. BACKGROUND7

follows a multivariate normal distribution with mean $\text{vec}(\boldsymbol{\mu})$ and covariance $\boldsymbol{\Sigma} \equiv \boldsymbol{\Sigma}_M \otimes \cdots \otimes \boldsymbol{\Sigma}_1$, where \otimes represents the Kronecker product; and, for $m = 1, \dots, M$, $\mathbf{X}_{(m)}$ follows a matrix normal distribution (Gupta & Nagar 2018) with mean $\boldsymbol{\mu}_{(m)}$, row covariance $\boldsymbol{\Sigma}_m$, and column covariance $\boldsymbol{\Sigma}_{-m} \equiv \boldsymbol{\Sigma}_M \otimes \cdots \otimes \boldsymbol{\Sigma}_{m+1} \otimes \boldsymbol{\Sigma}_{m-1} \otimes \cdots \otimes \boldsymbol{\Sigma}_1$.

2.2 Tensor Envelope

A tensor envelope (Li & Zhang 2017, Zhang & Li 2017) is a generalization of the envelopes in multivariate analysis (Cook et al. 2010, Cook 2018) that combines the Tucker tensor decomposition with the notion of reducing subspaces in functional analysis (Conway 2013). We briefly review these concepts below.

Given a matrix $\mathbf{B} \in \mathbb{R}^{p \times d}$, $\mathcal{B} = \text{span}(\mathbf{B}) \subseteq \mathbb{R}^p$ is defined as the subspace spanned by the column vectors of \mathbf{B} . Projections onto \mathcal{B} and its orthogonal complement subspace \mathcal{B}^\perp are denoted as $\mathbf{P}_\mathbf{B} = \mathbf{P}_\mathcal{B}$ and $\mathbf{Q}_\mathbf{B} = \mathbf{Q}_\mathcal{B} = \mathbf{I}_p - \mathbf{P}_\mathbf{B}$, respectively. If the matrix is of full column rank, then $\mathbf{P}_\mathbf{B} = \mathbf{B}(\mathbf{B}^T \mathbf{B})^{-1} \mathbf{B}$.

Definition 1. A subspace $\mathcal{B} \subseteq \mathbb{R}^p$ is a reducing subspace of $\mathbf{M} \in \mathbb{R}^{p \times p}$ if and only if $\mathbf{M} = \mathbf{P}_\mathbf{B} \mathbf{M} \mathbf{P}_\mathbf{B} + \mathbf{Q}_\mathbf{B} \mathbf{M} \mathbf{Q}_\mathbf{B}$. The \mathbf{M} -envelope of \mathcal{B} is the intersection of all reducing subspaces of \mathbf{M} that contain \mathcal{B} , and is denoted as $\mathcal{E}_\mathbf{M}(\mathcal{B})$ or $\mathcal{E}_\mathbf{M}(\mathbf{B})$.

By construction, the envelope $\mathcal{E}_\mathbf{M}(\mathbf{B}) \subseteq \mathbb{R}^p$ is always unique and the smallest such subspace. The existence is easily guaranteed when $\mathbf{M} > 0$ (Cook et al. 2010). Cook & Zhang (2015) provide a general envelope construction for a wide range of multivariate parameter estimation problems. In the gen-

2. BACKGROUND8

eral envelope construction, \mathbf{B} is the parameter of interest, and \mathbf{M} is either the covariance matrix of some random vector or the asymptotic covariance matrix of a \sqrt{n} -consistent estimator. Thus, the existence and uniqueness of envelopes are always true. Tensor envelopes have been developed under linear regression models with a tensor response (Li & Zhang 2017) and with a tensor predictor (Zhang & Li 2017). We unify the two envelopes, and give the following more general formulation of a tensor envelope. Let $\mathbf{B} \in \mathbb{R}^{p_1 \times \cdots \times p_M}$ be the tensorial parameter of interest. Let $\Sigma = \bigotimes_{m=1}^M \Sigma_m \equiv \Sigma_M \otimes \cdots \otimes \Sigma_1$ be the Kronecker product of a series of symmetric positive-definite matrices $\Sigma_m \in \mathbb{R}^{p_m \times p_m}$, for $m = 1, \dots, M$. The Kronecker operator for two subspaces $\mathcal{S}_1 \otimes \mathcal{S}_2$ is defined as the subspace spanned by $\mathbf{B}_1 \otimes \mathbf{B}_2$, where \mathbf{B}_k is any basis matrix for subspace \mathcal{S}_k , for $k = 1, 2$. The definition and a key property of the tensor envelope are summarized as follows.

Definition 2. *The tensor envelope $\mathcal{T}_\Sigma(\mathbf{B})$ is the intersection of all reducing subspaces \mathcal{S} of Σ that contain $\text{vec}(\mathbf{B})$ and can be written as $\mathcal{S} = \mathcal{S}_M \otimes \cdots \otimes \mathcal{S}_1$, with $\mathcal{S}_m \subseteq \mathbb{R}^{p_m}$, for $m = 1, \dots, M$.*

Proposition 1. $\mathcal{T}_\Sigma(\mathbf{B}) = \mathcal{E}_{\Sigma_M}(\mathbf{B}_{(M)}) \otimes \cdots \otimes \mathcal{E}_{\Sigma_1}(\mathbf{B}_{(1)})$.

This proposition (derived from Li & Zhang 2017, Proposition 3) connects the tensor envelope $\mathcal{T}_\Sigma(\mathbf{B})$ with the multivariate envelopes $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$, for $m = 1, \dots, M$, along each mode of the tensor \mathbf{B} . It implies that we can estimate a tensor envelope by estimating the individual envelopes $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$, for each mode m . The existence, uniqueness, and minimality of the envelope $\mathcal{E}_{\Sigma_M}(\mathbf{B}_{(m)})$

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE

are shown in Cook et al. (2010). Then, by Proposition 1, the tensor envelope $\mathcal{T}_{\Sigma}(\mathbf{B})$ always exists and is unique.

3. Discriminant Analysis with Tensor Envelope

3.1 The TDA Model

We consider an M th-order tensor variable $\mathbf{X} \in \mathbb{R}^{p_1 \times \dots \times p_M}$, for $M \geq 2$, and a categorical response $Y \in \{1, \dots, K\}$, with $K \geq 2$ categories/classes. We consider the following tensor discriminant analysis (TDA) model that is a natural generalization of the linear discriminant analysis model in the vector case:

$$\mathbf{X} \mid (Y = k) \sim TN(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_M), \quad (3.1)$$

where $\boldsymbol{\mu}_k \in \mathbb{R}^{p_1 \times \dots \times p_M}$ is the class-specific mean, and $\boldsymbol{\Sigma}_m \in \mathbb{R}^{p_m \times p_m}$, for $m = 1, \dots, M$, are symmetric positive-definite common covariance structures across classes. We assume nontrivial classes such that $\Pr(Y = k) = \pi_k > 0$ and $\sum_{k=1}^K \pi_k = 1$. Pan et al. (2019) also considered this model and developed a sparse TDA method. For a discriminant analysis, our goal is to improve the estimation of the optimal prediction of Y , which is known as Bayes' rule. Under the TDA model (3.1), the Bayes rule is given as

$$\phi^{\text{TDA}}(\mathbf{X}) = \operatorname{argmax}_{k=1, \dots, K} \Pr(Y = k \mid \mathbf{X}) = \operatorname{argmax}_{k=1, \dots, K} \{c_k + \langle \mathbf{B}_k, \mathbf{X} \rangle\}, \quad (3.2)$$

where $\mathbf{B}_k = \llbracket \boldsymbol{\mu}_k - \boldsymbol{\mu}_1; \boldsymbol{\Sigma}_1^{-1}, \dots, \boldsymbol{\Sigma}_M^{-1} \rrbracket$, $c_k = \log(\pi_k/\pi_1) - \langle \mathbf{B}_k, \frac{\boldsymbol{\mu}_k + \boldsymbol{\mu}_1}{2} \rangle$, and $\langle \mathbf{B}_k, \mathbf{X} \rangle$ is the inner product of \mathbf{B}_k and \mathbf{X} . Let $\mathbf{B} \in \mathbb{R}^{p_1 \times \dots \times p_M \times (K-1)}$ be the stacked tensor coefficients $\{\mathbf{B}_2, \dots, \mathbf{B}_K\}$. Then, $\phi^{\text{TDA}}(\mathbf{X})$ can be viewed as a

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹⁰

function of $\mathbf{B}_{(M+1)} \text{vec}(\mathbf{X}) = (\langle \mathbf{B}_2, \mathbf{X} \rangle, \dots, \langle \mathbf{B}_K, \mathbf{X} \rangle)^T \in \mathbb{R}^{K-1}$. Therefore, to improve classification accuracy, we need to improve the estimation of the tensor parameter \mathbf{B} . Although the TDA model reduces the number of parameters in the covariance matrix, many model parameters remain in most tensor data sets. Therefore, we use the tensor envelope to further reduce the number of parameters in the TDA model, thus facilitating the estimation.

3.2 The DATE Model

Similarly to the motivation of the envelope discriminant analysis for a vector predictor (Zhang & Mai 2019), the tensor envelope for a discriminant analysis and classification aims to identify and eliminate the part of $\mathbf{X} \in \mathbb{R}^{p_1 \times \dots \times p_M}$ that is unrelated to Bayes' classification rule and the remaining information. We consider a decomposition in the form of $\mathbf{X} = \mathbb{P}(\mathbf{X}) + \mathbb{Q}(\mathbf{X})$, where $\mathbb{P}(\mathbf{X}) = \llbracket \mathbf{X}; \mathbf{P}_1, \dots, \mathbf{P}_M \rrbracket$, and each $\mathbf{P}_m \in \mathbb{R}^{p_m \times p_m}$ is the projection onto a latent subspace $\mathcal{S}_m \subseteq \mathbb{R}^{p_m}$. The Tucker product form of $\mathbb{P}(\mathbf{X})$ preserves all information for discriminating Y , and $\mathbb{Q}(\mathbf{X}) = \mathbf{X} - \mathbb{P}(\mathbf{X})$ is the part that is irrelevant for classification. As such, we consider that, for $k = 1, \dots, K$,

$$\Pr(Y = k | \mathbf{X}) = \Pr\{Y = k | \mathbb{P}(\mathbf{X})\}, \quad \mathbb{Q}(\mathbf{X}) \perp \mathbb{P}(\mathbf{X}) | (Y = k), \quad (3.3)$$

where \perp indicates the independence of random variables. The first condition in (3.3) implies that Bayes' classification rule does not change if we project the data onto the subspaces \mathcal{S}_m , for $m = 1, \dots, M$, along each mode of the tensor. The second condition in (3.3) requires that the material part $\mathbb{P}(\mathbf{X})$ is not affected

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹¹

by the immaterial part $\mathbb{Q}(\mathbf{X})$. The following proposition connects the subspaces \mathcal{S}_m , for $m = 1, \dots, M$, with the TDA model parameters and Bayes' rule.

Proposition 2. *Under the TDA model (3.1), assumption (3.3) is equivalent to*

$$\text{span}(\mathbf{B}_{(M+1)}^T) \subseteq \mathcal{S} = \mathcal{S}_M \otimes \dots \otimes \mathcal{S}_1, \quad \Sigma_m = \mathbf{P}_m \Sigma_m \mathbf{P}_m + \mathbf{Q}_m \Sigma_m \mathbf{Q}_m,$$

for $m = 1, \dots, M$. Furthermore, $\mathcal{E}_{\Sigma_{\mathbf{X}}}(\mathbf{B}_{(M+1)}^T) \subseteq \mathcal{T}_{\Sigma}(\mathbf{B}) = \mathcal{E}_{\Sigma_M}(\mathbf{B}_{(M)}) \otimes \dots \otimes \mathcal{E}_{\Sigma_1}(\mathbf{B}_{(1)})$, where $\Sigma_{\mathbf{X}} = \text{cov}(\text{vec}(\mathbf{X}))$.

Proposition 2 establishes the tensor envelope construct $\mathcal{T}_{\Sigma}(\mathbf{B})$, and implies that $\Pr(Y = k | \mathbf{X}) = P(Y = k | \mathbf{X} \times_m \mathbf{P}_m)$ and $\mathbf{X} \times_m \mathbf{P}_m \perp \mathbf{X} \times_m \mathbf{Q}_m | (Y = k)$ on each mode $m = 1, \dots, M$, where $\mathbf{Q}_m = \mathbf{I}_{p_m} - \mathbf{P}_m$ is the projection onto \mathcal{S}_m^{\perp} . When $M = 1$, this reduces to the envelope LDA model (Zhang & Mai 2019). Proposition 2 offers a more intuitive explanation of the tensor envelope: it is the smallest subspace reduction $\mathbb{P}(\mathbf{X})$ to attain the same Bayes' rule as the original \mathbf{X} , while not leaking information by correlating with $\mathbb{Q}(\mathbf{X})$. Finally, the tensor envelope is shown to contain the vectorized envelope discriminant subspace $\mathcal{E}_{\Sigma_{\mathbf{X}}}(\mathbf{B}_{(M+1)}^T)$, and is a reducing subspace for the marginal covariance of $\text{vec}(\mathbf{X})$.

To further investigate the subspace representations in Proposition 2, we let $(\mathbf{\Gamma}_m, \mathbf{\Gamma}_{0m})$ be an orthogonal basis matrix for \mathbb{R}^{p_m} such that $\text{span}(\mathbf{\Gamma}_m) = \mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$. Let the envelope dimension be u_m , $u_m \leq p_m$, $\mathbf{\Gamma}_m \in \mathbb{R}^{p_m \times u_m}$, and $\mathbf{\Gamma}_{0m} \in \mathbb{R}^{p_m \times (p_m - u_m)}$. Then, we have the following parameterization of the TDA

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹²

model, for $k = 1, \dots, K$ and $m = 1, \dots, M$:

$$\mathbf{B}_k = \llbracket \boldsymbol{\eta}_k; \boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M \rrbracket, \boldsymbol{\eta}_k \in \mathbb{R}^{u_1 \times \dots \times u_M}, \quad (3.4)$$

$$\boldsymbol{\Sigma}_m = \boldsymbol{\Gamma}_m \boldsymbol{\Omega}_m \boldsymbol{\Gamma}_m^T + \boldsymbol{\Gamma}_{0m} \boldsymbol{\Omega}_{0m} \boldsymbol{\Gamma}_{0m}^T, \quad (3.5)$$

where $\boldsymbol{\eta}_k$, defined as a tensor with conforming dimensions, consists of the coordinates of \mathbf{B}_k with respect to the basis matrices $\boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M$. Similarly, $\boldsymbol{\Omega}_m, \boldsymbol{\Omega}_{0m}$ and the following $\boldsymbol{\Theta}_k$ are constructed based on the basis matrices $\boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M$. Note that $\bar{\boldsymbol{\mu}} = \sum_{k=1}^K \pi_k \boldsymbol{\mu}_k$ is the overall mean tensor. It is straightforward to show that (3.4) is equivalent to the following equation:

$$\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}} = \llbracket \boldsymbol{\Theta}_k; \boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M \rrbracket, \boldsymbol{\Theta}_k \in \mathbb{R}^{u_1 \times \dots \times u_M}, \quad (3.6)$$

where we have an additional constraint $\sum_{k=1}^K \pi_k \boldsymbol{\Theta}_k = 0$. The total number of free parameters in the TDA model (3.1) is $(K-1) + K \prod_{m=1}^M p_m + \sum_{m=1}^M \frac{p_m(p_m+1)}{2}$, whereas the total number of free parameters in (3.5) and (3.6) is $(K-1) + K \prod_{m=1}^M u_m + \sum_{m=1}^M \frac{p_m(p_m+1)}{2}$. Therefore, the tensor envelope has reduced the number of parameters by $K(\prod_{m=1}^M p_m - \prod_{m=1}^M u_m)$ under the TDA model. By reducing the model complexity, the envelope approach can often lead to a substantial gain in the estimation of \mathbf{B}_k , thus improving the classification.

3.3 Likelihood-based Estimation

Given the observed data $\{\mathbf{X}^i, Y^i\}_{i=1}^n$, for each class $k = 1, \dots, K$, we have $n_k = \sum_i I(Y^i = k)$ as the class- k sample size and $\bar{\mathbf{X}}_k = n_k^{-1} \sum_i I(Y^i = k) \mathbf{X}^i$ as the class- k sample mean. Under TDA model (3.1), the standard maximum

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹³

likelihood estimators (MLEs) for π_k and $\boldsymbol{\mu}_k$ are n_k/n and $\bar{\mathbf{X}}_k$, respectively. In addition, the MLE for $\boldsymbol{\Sigma}_m$ can be obtained through iterative updates as the solution to

$$\boldsymbol{\Sigma}_m = \frac{1}{np_{-m}} \sum_{i=1}^n I(Y^i = K) (\mathbf{X}^i - \bar{\mathbf{X}}_k)_{(m)} \boldsymbol{\Sigma}_{-m}^{-1} (\mathbf{X}^i - \bar{\mathbf{X}}_k)_{(m)}^T, \quad (3.7)$$

where $p_{-m} = \prod_{l \neq m} p_l$, and $\boldsymbol{\Sigma}_{-m} = \boldsymbol{\Sigma}_M \otimes \cdots \otimes \boldsymbol{\Sigma}_{m+1} \otimes \boldsymbol{\Sigma}_{m-1} \otimes \cdots \otimes \boldsymbol{\Sigma}_1$. The derivation for (3.7) is similar to that in Manceur & Dutilleul (2013), and is thus omitted. Under the tensor envelope parameterization (3.4) and (3.5), we have derived MLE equations to greatly facilitate estimation.

Proposition 3. *Under the DATE model (3.1), (3.4), and (3.5), the MLE of the envelope basis $\hat{\boldsymbol{\Gamma}}_m$ is obtained by minimizing the following objective function under the semi-orthogonal constraint $\boldsymbol{\Gamma}_m^T \boldsymbol{\Gamma}_m = \mathbf{I}_{p_m}$:*

$$\ell_m(\boldsymbol{\Gamma}_m) = \log |\boldsymbol{\Gamma}_m^T \hat{\mathbf{M}}_m \boldsymbol{\Gamma}_m| + \log |\boldsymbol{\Gamma}_m^T (\hat{\mathbf{N}}_m)^{-1} \boldsymbol{\Gamma}_m|, \quad (3.8)$$

where $\hat{\mathbf{M}}_m = (n \prod_{m \neq j} p_m)^{-1} \sum_{i=1}^n I(Y^i = K) \{\mathbf{s}_{(m)}^{ik} \hat{\boldsymbol{\Sigma}}_{-m}^{-1} (\mathbf{s}_{(m)}^{ik})^T\}$, $\hat{\mathbf{N}}_m = (n \prod_{m \neq j} p_m)^{-1} \sum_{i=1}^n (\mathbf{X}^i - \bar{\mathbf{X}})_{(m)} \hat{\boldsymbol{\Sigma}}_{-m}^{-1} (\mathbf{X}^i - \bar{\mathbf{X}})_{(m)}^T$, and $\mathbf{s}^{ik} = I(Y^i = k) (\mathbf{X}^i - \bar{\mathbf{X}}) - \llbracket \bar{\mathbf{X}}_k - \bar{\mathbf{X}}; \hat{\mathbf{P}}_1, \dots, \hat{\mathbf{P}}_{m-1}, \mathbf{I}_{p_m}, \hat{\mathbf{P}}_{m+1}, \dots, \hat{\mathbf{P}}_M \rrbracket$. The MLEs for the DATE parameters are given by

$$\begin{aligned} \hat{\boldsymbol{\Theta}}_k &= \llbracket \bar{\mathbf{X}}_k - \bar{\mathbf{X}}; \hat{\boldsymbol{\Gamma}}_1^T, \dots, \hat{\boldsymbol{\Gamma}}_M^T \rrbracket, \quad \hat{\mathbf{B}}_k = \llbracket \bar{\mathbf{X}}_k - \bar{\mathbf{X}}_1; \hat{\boldsymbol{\Gamma}}_1 \hat{\boldsymbol{\Omega}}_1^{-1} \hat{\boldsymbol{\Gamma}}_1^T, \dots, \hat{\boldsymbol{\Gamma}}_M \hat{\boldsymbol{\Omega}}_M^{-1} \hat{\boldsymbol{\Gamma}}_M^T \rrbracket, \\ \hat{\boldsymbol{\Sigma}}_m &= \hat{\boldsymbol{\Gamma}}_m \hat{\boldsymbol{\Omega}}_m \hat{\boldsymbol{\Gamma}}_m^T + \hat{\boldsymbol{\Gamma}}_{0m} \hat{\boldsymbol{\Omega}}_{0m} \hat{\boldsymbol{\Gamma}}_{0m}^T, \quad \hat{\boldsymbol{\Omega}}_m = \frac{1}{np_{-m}} \sum_{i=1}^n I(Y^i = K) \{\hat{\boldsymbol{\Gamma}}_m^T \mathbf{s}_{(m)}^{ik} \hat{\boldsymbol{\Sigma}}_{-m}^{-1} (\mathbf{s}_{(m)}^{ik})^T \hat{\boldsymbol{\Gamma}}_m\}, \\ \hat{\boldsymbol{\Omega}}_{0m} &= \frac{1}{np_{-m}} \sum_{i=1}^n \{\hat{\boldsymbol{\Gamma}}_{0m}^T (\mathbf{X}^i - \bar{\mathbf{X}})_{(m)} \hat{\boldsymbol{\Sigma}}_{-m}^{-1} (\mathbf{X}^i - \bar{\mathbf{X}})_{(m)}^T \hat{\boldsymbol{\Gamma}}_{0m}\}. \end{aligned}$$

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹⁴

The implementation of the algorithm is based directly on the above proposition. To initialize, we first obtain the standard MLE for Σ_m based on equation (3.7), and replace $\hat{\Gamma}_m$ with \mathbf{I}_{p_m} in the construction of the pseudo observations \mathbf{s}^{ik} . The pseudo observations are used to update the envelope basis $\hat{\Gamma}_m$ by minimizing $\ell_m(\Gamma_m)$ in Proposition 3. We then iteratively update the envelope basis, the pseudo observations, and the covariance matrices $\hat{\Sigma}_m$. After convergence, we calculate the MLEs for the means μ_k and for other parameters, such as Θ_k and \mathbf{B} . Finally, after we obtain $\hat{\mathbf{B}}$, the prediction is simply the LDA classification rule (3.2) on the reduced data $(\langle \hat{\mathbf{B}}_2, \mathbf{X} \rangle, \dots, \langle \hat{\mathbf{B}}_K, \mathbf{X} \rangle)^T \in \mathbb{R}^{K-1}$.

The objective function $\ell_m(\Gamma_m)$ in (3.8) is nonconvex and depends on all other $\{\Gamma_j\}_{j \neq m}$ intrinsically through $\hat{\Sigma}_{-m}$ and \mathbf{s}^{ik} . Therefore, the fully iterative algorithm can be slow and sensitive to initialization. To speed up the computation, we adopt the one-step estimation procedure of Li & Zhang (2017) for tensor envelopes. Specifically, we run the iteration of the MLE equations only once. That is, without alternately updating all Γ_m , we optimize each $\ell_m(\Gamma_m)$ separately, which is sped up further by the 1D envelope algorithm (Cook & Zhang 2016). Because the one-step estimation in the DATE model is similar to that in Li & Zhang (2017), the implementation details are relegated to Section S2 of the Supplementary Material. In practice, note that the estimation accuracy (e.g., for the key parameter \mathbf{B}) and the classification/prediction accuracy of the one-step estimator are always as good as the MLEs. Such findings are consistent with those of previous works (Cook & Zhang 2016, Li & Zhang 2017), where the one-step estimator and the 1D algorithm are shown to outperform to be better

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹⁵

than the full MLE updates in practice. Therefore, we present only the numerical results obtained using the one-step estimator in our simulations and real-data analysis.

Another important hyperparameter that we need to estimate is the envelope dimension. To select the envelope dimension u_m , for $m = 1, \dots, M$, we apply cross-validation with a grid search to choose u_m that minimizes the classification error. Because cross-validation tends to overfit the envelope dimension, we adopt the “one standard error rule”, in which we choose the smallest u_m with an error that is no more than one standard error above the minimum cross-validated error. This method has proven to be stable in simulation studies. Note that the Bayesian information criterion (BIC) is frequently considered for envelope dimension selection. Owing to the carefully derived and simplified objective function $\ell_m(\Gamma_m)$ in (3.8), we may directly apply the 1D-BIC envelope dimension selection (Zhang & Mai 2018), which has been proven theoretically and is computationally feasible (it is much faster and more stable than the standard BIC in envelope dimension selection). See Section S2.2 of the Supplementary Material for additional discussion and numerical examples.

3.4 Decomposition-based Estimation

Our decomposition-based approach is motivated by the following lemma.

Lemma 1. *Under parameterization (3.4) and (3.5), $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)}) = \mathcal{E}_{\Sigma_m}(\mathbf{U}_m)$, for $m = 1, \dots, M$, where $\mathbf{U}_m = p_{-m}^{-1} \{ \sum_{k=1}^K \pi_k (\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}})_{(m)} (\boldsymbol{\mu}_k - \bar{\boldsymbol{\mu}})_{(m)}^T \}$. If we further assume that the eigenvalues of $\mathbf{P}_{\Gamma_m} \Sigma_m \mathbf{P}_{\Gamma_m}$ are distinct from those*

3. DISCRIMINANT ANALYSIS WITH TENSOR ENVELOPE¹⁶

of $\mathbf{Q}_{\Gamma_m} \Sigma_m \mathbf{Q}_{\Gamma_m}$, then $\mathcal{E}_{\Sigma_m}(\mathbf{U}_m) = \sum_{(\mathbf{v}_i^{(m)})^T \mathbf{U}_m \mathbf{v}_i^{(m)} \neq 0} \text{span}(\mathbf{v}_i^{(m)})$, for $m = 1, \dots, M$, where $\mathbf{v}_i^{(m)}$ is the i th eigenvector of Σ_m .

Lemma 1 establishes the equivalence between $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$ and $\mathcal{E}_{\Sigma_m}(\mathbf{U}_m)$, where \mathbf{U}_m is a positive semi-definite symmetric $p_m \times p_m$ matrix. This matrix \mathbf{U}_m does not involve the covariance inverse Σ_m^{-1} , as in \mathbf{B} , and can be viewed as the mode- m between the class variance of \mathbf{X} . The symmetry of \mathbf{U}_m also facilitates the estimation later. In order to construct a decomposition-based method, we assume that the eigenvalues of $\mathbf{P}_{\Gamma_m} \Sigma_m \mathbf{P}_{\Gamma_m}$ are distinct from those of $\mathbf{Q}_{\Gamma_m} \Sigma_m \mathbf{Q}_{\Gamma_m}$. This assumption is mild, and much weaker than requiring Σ_m to have p_m distinct eigenvalues. Under this mild assumption, the tensor envelope can be obtained by recognizing $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$ as the subspace spanned by all eigenvectors $\mathbf{v}_i^{(m)}$ of Σ_m that are not orthogonal to \mathbf{U}_m , namely, $\sum_{(\mathbf{v}_i^{(m)})^T \mathbf{U}_m \mathbf{v}_i^{(m)} \neq 0} \text{span}(\mathbf{v}_i^{(m)})$.

Thus, the DATE-D procedure for $\mathcal{E}_{\Sigma_m}(\mathbf{B}_{(m)})$ in the population is as follows:

1. Obtain the eigenvectors of Σ_m : $\mathbf{v}_1^{(m)}, \dots, \mathbf{v}_{p_m}^{(m)}$, with ordered eigenvalues $\lambda_1^{(m)} \geq \dots \geq \lambda_{p_m}^{(m)}$.
2. Calculate the envelope scores: $\phi_l^{(m)} = (\mathbf{v}_l^{(m)})^T \mathbf{U}_m \mathbf{v}_l^{(m)}$, for $l = 1, \dots, p_m$.
3. Organize the envelope scores in descending order $\phi_{(1)}^{(m)} \geq \phi_{(2)}^{(m)} \geq \dots \geq \phi_{(p_m)}^{(m)}$, and let $\mathbf{v}_{(j)}^{(m)}$ be the eigenvector corresponding to $\phi_{(j)}^{(m)}$.
4. Output the envelope as $\mathcal{E}_{\Sigma_m}(\mathbf{U}_m) = \text{span}(\mathbf{v}_{(1)}^{(m)}, \dots, \mathbf{v}_{(u_m)}^{(m)})$, which has a basis matrix of $\Gamma_m = (\mathbf{v}_{(1)}^{(m)}, \dots, \mathbf{v}_{(u_m)}^{(m)})$.

The sample algorithm is readily available by replacing Σ_m and \mathbf{U}_m with their sample counterparts: $\tilde{\Sigma}_m = \frac{1}{np-m} \sum_{k=1}^K I(Y^i = K)(\mathbf{X}^i - \bar{\mathbf{X}}_k)_{(m)}(\mathbf{X}^i - \bar{\mathbf{X}}_k)_{(m)}^T$ and $\hat{\mathbf{U}}_m = \frac{1}{p-m} \left\{ \sum_{k=1}^K \frac{n_k}{n} (\bar{\mathbf{X}}_k - \bar{\mathbf{X}})_{(m)}(\bar{\mathbf{X}}_k - \bar{\mathbf{X}})_{(m)}^T \right\}$, respectively. Note that the estimate $\tilde{\Sigma}_m$ is a closed-form solution, which avoids the iterations within the covariance matrices in (3.7). This modification further accelerates the computation and facilitates theoretical analyses in high dimensions. The algorithm can be viewed as an extension of the algorithm in Zhang et al. (2021) for vector data.

The DATE-D procedure can be intuitively viewed as selecting the eigenvectors of Σ_m with nonzero envelope scores. The computationally most expensive part of DATE-D is the eigen-decomposition of Σ_m in Step 1. Because no matrix inversion is needed, DATE-D can be applied to very high-dimensional settings. Striving for a best prediction, we again use cross-validation to select the envelope dimension for DATE-D. Note that the selected “most predictive” envelope dimensions may differ for DATE-L and DATE-D.

4. Theory

4.1 Asymptotic Properties of DATE-L

Here, we study the consistency and asymptotic efficiency of three likelihood-based estimators: the MLE without the separable covariance assumption (i.e., the standard LDA on vectorized data); the MLE with the separable covariance assumption (i.e., the TDA); and the MLE under the DATE model (i.e., DATE-L). Though we recommend using the one-step estimation rather than the full MLE

updates for DATE-L in practice, the asymptotic properties of the MLE provide an idealistic “best-case scenario” and insights into the potential advantages of the DATE model. The \sqrt{n} -consistency of the one-step estimator is provided in the Supplementary Material (Theorem S1).

The results are presented for all the parameters in the model, with additional focus on the estimation of \mathbf{B}_k , for $k = 2, \dots, K$, or equivalently, $\beta_k \equiv \text{vec}(\mathbf{B}_k)$. The estimators are denoted by $\hat{\beta}_k^{\text{LDA}} = \mathbf{S}^{-1} \text{vec}(\bar{\mathbf{X}}_k - \bar{\mathbf{X}}_1)$ for the LDA estimator, $\hat{\beta}_k^{\text{TDA}} = (\hat{\Sigma}_M \otimes \dots \otimes \hat{\Sigma}_1) \text{vec}(\bar{\mathbf{X}}_k - \bar{\mathbf{X}}_1)$ for the TDA estimator, and $\hat{\beta}_k^{\text{DATE}} = (\hat{\Gamma}_M \otimes \dots \otimes \hat{\Gamma}_1) \text{vec}(\hat{\eta}_k)$ for DATE-L.

We first compare the asymptotic efficiency of $\hat{\beta}_k^{\text{LDA}}$, $\hat{\beta}_k^{\text{TDA}}$, and $\hat{\beta}_k^{\text{DATE}}$. Specifically, we define the parameter vectors corresponding to each estimator as follows, where we stack all unique parameters in a vector using the operators vec (vectorization of matrix/tensor) and vech (vectorization of symmetric matrix by stacking the lower triangular of the matrix), $\mathbf{h}^T = (\{\beta_k^T\}_{k=2}^K, \text{vech}^T(\Sigma)), \boldsymbol{\eta}^T = (\{\beta_k^T\}_{k=2}^K, \{\text{vech}^T(\Sigma_m)\}_{m=1}^M)$, and for envelope parameters $\boldsymbol{\xi} = (\{\text{vec}^T(\Gamma_m)\}_{m=1}^M, \{\text{vec}^T(\boldsymbol{\eta}_k)\}_{k=2}^K, \{\text{vech}^T(\Omega_m)\}_{m=1}^M, \{\text{vech}^T(\Omega_{0m})\}_{m=1}^M)$. It is straightforward to calculate the number of parameters based on the length of each parameter vector. Specifically, \mathbf{h} corresponds to the parameterization in the vectorized LDA model, where the covariance $\Sigma \in \mathbb{R}^{\prod_m p_m \times \prod_m p_m}$ is unstructured; $\boldsymbol{\eta}$ corresponds to the parameterization in the TDA model with the separable Kronecker covariance structure; finally, $\boldsymbol{\xi}$ contains all parameters under the tensor envelope structure.

From (3.4) and (3.5), we can see that \mathbf{h} is an estimable function of $\boldsymbol{\eta}$ and $\boldsymbol{\xi}$;

that is, $\mathbf{h} = \mathbf{h}(\boldsymbol{\eta}) = \mathbf{h}(\boldsymbol{\xi})$. We define the gradient as $\mathbf{H} = \frac{\partial \mathbf{h}(\boldsymbol{\eta})}{\partial \boldsymbol{\phi}}$ and $\mathbf{K} = \frac{\partial \mathbf{h}(\boldsymbol{\xi})}{\partial \boldsymbol{\xi}}$. We denote $\widehat{\mathbf{h}}_{\text{TDA}}$ as the standard MLEs containing the sample estimators $\overline{\mathbf{X}}_k$ and $\widehat{\Sigma}_m$ in (3.7). Similarly, $\widehat{\mathbf{h}}_{\text{DATE}}$ and $\widehat{\mathbf{h}}_{\text{LDA}}$ are the MLEs under the DATE model and the vectorized LDA model, respectively.

Theorem 1. *Assume (\mathbf{X}^i, Y^i) , for $i = 1, \dots, n$, are independent and identically distributed (i.i.d) according to the DATE model (3.1), (3.4), and (3.5). Then $\sqrt{n} \text{vec}(\widehat{\boldsymbol{\beta}}^{\text{LDA}} - \boldsymbol{\beta}) \rightarrow_d N(0, \mathbf{W}_\beta)$; $\sqrt{n} \text{vec}(\widehat{\boldsymbol{\beta}}^{\text{TDA}} - \boldsymbol{\beta}) \rightarrow_d N(0, \mathbf{U}_\beta)$; and $\sqrt{n} \text{vec}(\widehat{\boldsymbol{\beta}}^{\text{DATE}} - \boldsymbol{\beta}) \rightarrow_d N(0, \mathbf{V}_\beta)$. Moreover, $\mathbf{V}_\beta \leq \mathbf{U}_\beta \leq \mathbf{W}_\beta$.*

The detailed expressions of the asymptotic variances \mathbf{V}_β , \mathbf{U}_β , and \mathbf{W}_β are provided in the Supplementary Material (Section S4.2). Theorem 1 establishes the \sqrt{n} -consistency and asymptotic normality of all three types of MLEs. The result is not surprising, because we gain more efficiency by using more structures, while maximizing the likelihood.

To gain further insights, we consider the oracle envelope estimator of $\boldsymbol{\beta}$, denoted as $\widehat{\boldsymbol{\beta}}_\Gamma$, that replaces $\widehat{\Gamma}$ with the true envelope basis Γ in the estimation.

Theorem 2. *Under the same conditions as those in Theorem 1, $\widehat{\boldsymbol{\beta}}_\Gamma$ is \sqrt{n} -consistent and asymptotically normal. The asymptotic covariance of $\text{vec}(\widehat{\boldsymbol{\beta}}_\Gamma)$ is $\mathbf{V}_\Gamma = \mathbf{A} \otimes \{(\Gamma_M \Omega_M^{-1} \Gamma_M^T) \otimes \dots \otimes (\Gamma_1 \Omega_1^{-1} \Gamma_1^T)\}$, where $\mathbf{A} = \text{diag}(\pi_2^{-1}, \dots, \pi_K^{-1}) + \pi_1^{-1} \mathbf{1}_{K-1} \mathbf{1}_{K-1}^T$ is a constant matrix.*

Because we can write $\Sigma = (\Gamma_M \Omega_M \Gamma_M^T + \Gamma_{0M} \Omega_{0M} \Gamma_{0M}^T) \otimes \dots \otimes (\Gamma_1 \Omega_1 \Gamma_1^T + \Gamma_{01} \Omega_{01} \Gamma_{01}^T)$, $\mathbf{V}_\Gamma \leq \mathbf{U}_\beta$. In particular, a direct comparison shows that the envelope estimators have bigger potential gains in efficiency when the immaterial

variation Ω_{0m} is large relative to the material variation Ω_m in the predictor.

Finally, we establish the \sqrt{n} -consistency and asymptotic normality of the envelope estimator under mild moment conditions, instead of the tensor normality assumption in (3.1). Specifically, we consider the tensor envelope parameterization without the distributional assumption. Then, the conditional independence assumption in (3.3), $\mathbb{Q}(\mathbf{X}) \perp \mathbb{P}(\mathbf{X}) \mid (Y = k)$, is weakened to the assumption of uncorrelated $\mathbb{Q}(\mathbf{X})$ and $\mathbb{P}(\mathbf{X})$ in each class $Y = k$.

Theorem 3. *For $k = 1, \dots, K$, assume $\text{vec}(\mathbf{X}^i) \mid Y^i = k$, for $i = 1, \dots, n$ are i.i.d with finite fourth moments with mean tensor $\boldsymbol{\mu}_k$ and separable covariance $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_M \otimes \dots \otimes \boldsymbol{\Sigma}_1$, and parameterizations (3.4) and (3.5) are still satisfied. Then $\sqrt{n}\text{vec}(\widehat{\mathbf{h}}_{\text{TDA}})$ converges to a normal distributions with mean zero and asymptotic covariance matrix $\text{avar}(\sqrt{n}\widehat{\mathbf{h}}_{\text{TDA}}) = \mathbf{H}(\mathbf{H}^T \mathbf{J}_h \mathbf{H})^\dagger \mathbf{H}^T \mathbf{J}_h \boldsymbol{\Xi} \mathbf{J}_h \mathbf{H}(\mathbf{H}^T \mathbf{J}_h \mathbf{H})^\dagger \mathbf{H}^T$, and $\sqrt{n}\text{vec}(\widehat{\mathbf{h}}_{\text{DATE}})$ converges to a normal distribution with mean zero and asymptotic covariance matrix $\text{avar}(\sqrt{n}\widehat{\mathbf{h}}_{\text{DATE}}) = \mathbf{K}(\mathbf{K}^T \mathbf{J}_h \mathbf{K})^\dagger \mathbf{K}^T \mathbf{J}_h \boldsymbol{\Xi} \mathbf{J}_h \mathbf{K}(\mathbf{K}^T \mathbf{J}_h \mathbf{K})^\dagger \mathbf{K}^T$. Furthermore, $\text{avar}(\sqrt{n}\widehat{\mathbf{h}}_{\text{DATE}}) \leq \text{avar}(\sqrt{n}\widehat{\mathbf{h}}_{\text{TDA}}) \leq \text{avar}(\sqrt{n}\widehat{\mathbf{h}}^{\text{LDA}})$ if $\text{span}(\mathbf{J}_h^{1/2} \mathbf{H})$ and $\text{span}(\mathbf{J}_h^{1/2} \mathbf{K})$ are reducing subspaces of $\mathbf{J}_h^{1/2} \boldsymbol{\Xi} \mathbf{J}_h^{1/2}$, where $\boldsymbol{\Xi} = \text{avar}(\sqrt{n}\widehat{\mathbf{h}}^{\text{LDA}})$.*

Theorem 3 shows that the envelope estimator is robust to model misspecification in the sense that it is \sqrt{n} -consistent without tensor normality. Moreover, the DATE-L estimator still has potential advantages over the standard TDA estimator for nonnormal data (See Section 5.3 for simulation examples).

The classification error rate obtained from the Bayes rule is a continuous function of the parameters $(\boldsymbol{\beta}_k, \pi_k, \boldsymbol{\mu}_k)$, for $k = 1, \dots, K$, under the LDA

model. Then, by the delta method, the asymptotic efficiency gain by DATE-L established in the above theorems implies a more accurate classification error. This is analogous to the efficiency gain and classification error rate comparison of the MLE versus the logistic regression (a \sqrt{n} -consistent, but asymptotically less efficient estimator) under the LDA model (Efron 1975, Bi & Jeske 2010).

4.2 Theoretical properties of DATE-D

We establish the convergence rate of the DATE-D estimator in high dimensions, where p_m can grow faster than n . We use c and C to represent generic positive constants that may vary. For simplicity, the envelope dimensions u_1, \dots, u_M are treated as constants that do not grow with p or n . We first introduce some technical assumptions:

- (A1) The eigenvalues of Σ_m , for $m = 1, \dots, M$, are all bounded between positive constants c_1 and c_2 .
- (A2) The smallest nonzero eigenvalue of U_m is bounded below by c_3 .
- (A3) $\|\mu_k - \mu\|_F \leq c_4$.
- (A4) The difference between any eigenvalue of $P_{\Gamma_m} \Sigma_m P_{\Gamma_m}$ and each $Q_{\Gamma_m} \Sigma_m Q_{\Gamma_m}$ is greater than c_5 .
- (A5) $c_6/K \leq n_k/n \leq c_7/K$, for $k = 1, \dots, K$.

Assumption (A1) implies that the population parameter Σ_m is wellconditioned, regardless of how p_m grows. Assumption (A2) can be viewed as a “signal

strength assumption” that ensures that the envelope scores calculated from \mathbf{U}_m are sufficiently accurate. Assumption (A3) is a mild assumption, because $\boldsymbol{\mu}_k - \boldsymbol{\mu} = \llbracket \boldsymbol{\theta}_k^*; \boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M \rrbracket$, for some $\boldsymbol{\theta}_k^* \in \mathbb{R}^{u_1 \times \dots \times u_m}$. Because $\boldsymbol{\theta}_k^*$ is a low-dimensional tensor, it is natural to assume that $\|\boldsymbol{\theta}_k^*\|_F \leq c_3$. Then, we arrive at Assumption (A3) by noting that $\|\boldsymbol{\mu}_k - \boldsymbol{\mu}\|_F = \|\boldsymbol{\theta}_k^*\|_F$. Assumption (A4) is required for the identifiability of the envelope from the decomposition perspective (Lemma 1). Assumption (A5) guarantees that each class has a decent sample size. All five assumptions are satisfied in our simulation examples (M1)–(M4) under covariance (C1) in Section 5.2.

We use $\eta_m = \sqrt{p_m/p_{-m}}$ to quantify the squareness of the matricization $\mathbf{X}_{(m)}$. For a tensor in which no mode’s dimension dominates all other modes combined, η_m is small. We define the classification error rate formally as $\widehat{R} = \Pr(\widehat{Y}(\widehat{\mathbf{B}}_k, \widehat{c}_k, k = 1, \dots, K) \neq Y)$, where $\widehat{Y}(\widehat{\mathbf{B}}_k, \widehat{c}_k, k = 1, \dots, K) = \operatorname{argmax}_{k=1, \dots, K} \{\widehat{c}_k + \langle \widehat{\mathbf{B}}_k, \mathbf{X} \rangle\}$. The population counterpart R is thus the Bayes error. Denote $\|\mathbf{A}\|_2$ for a tensor \mathbf{A} as the ℓ_2 -norm of $\operatorname{vec}(\mathbf{A})$.

Theorem 4. *Under assumptions (A1)–(A5), for a constant $C_3 > 1$, we have*

$$\|\mathbf{P}_{\widehat{\boldsymbol{\Gamma}}_m} - \mathbf{P}_{\boldsymbol{\Gamma}_m}\|_F = n^{-1/2}O(\eta_m + 1), \text{ for } m = 1, \dots, M,$$

with probability at least $1 - K \exp\{-C_1 p_m (C_3 - 1)\}$. Moreover,

$$\|\widehat{\mathbf{B}}_k - \mathbf{B}_k\|_2 = n^{-1/2}O(\max_m \eta_m + 1), \quad |\widehat{R} - R| = n^{-1/2}O(\max_m \eta_m + 1),$$

with probability at least $1 - C_2 K M \exp\{-C_1 p_m (C_3 - 1)\}$.

Corollary 1. *Under assumptions (A1)–(A5), when $n \gg \eta_m$, $p_m \rightarrow \infty$, and $n \rightarrow \infty$, we have $\mathbf{P}_{\widehat{\boldsymbol{\Gamma}}_m} \rightarrow \mathbf{P}_{\boldsymbol{\Gamma}_m}$, $\widehat{\mathbf{B}}_k \rightarrow \mathbf{B}_k$, and $\widehat{R} \rightarrow R$ in probability.*

The result in Theorem 4 is sufficiently strong for most tensor data applications, because p_{-m} is usually greater than p_m , especially when the $M \geq 3$. If the dimensions p_m , for $m = 1, \dots, M$, grow at the same rate, the ratio η_m either converges to zero ($M \geq 3$) or is bounded from above by a constant ($M = 2$). Then, we have \sqrt{n} -consistency for arbitrarily high-dimensional p_m when $M \geq 2$. However, for vector data, the rate becomes $(p/n)^{1/2}$, which means p cannot grow too fast. Hence, Theorem 4 reveals a fundamental difference between tensor and vector data. For vector data, it is challenging to estimate the covariance matrix accurately, but in tensor data, we can aggregate the information from different modes to achieve a consistent estimation of Σ_m .

5. Numerical Studies

5.1 Comparison setup

To investigate the empirical performance of the proposed DATE methods, we consider both simulations and real-data examples. In Section 5.2, we consider simulations under the DATE model. In Section 5.3, we consider models in which the DATE assumptions are violated. In Section 5.4, we construct a model with a high-dimensional matrix predictor to verify the consistency of DATE-D in high dimensions (cf., Theorem 4). In Section 5.5, we demonstrated our methods using real-data examples from colorimetric sensor arrays and longitudinal gene expressions.

We include various classification methods as competitors. First, we consider

the standard LDA and TDA estimators. However, owing to the high dimensionality, the standard LDA estimator is not applicable, and is hence replaced by the diagonal LDA (DLDA) (Dudoit et al. 2002). As regularized classification methods for high-dimensional vector data, we include an ℓ_1 -penalized Fisher's discriminant analysis (ℓ_1 -FDA; Witten & Tibshirani 2011) and an ℓ_1 -penalized logistic and multinomial logistic regression (ℓ_1 -GLM; Friedman et al. 2010). We also include several recent methods for matrix/tensor classification methods: a distance-weighted discrimination for multi-way data (m -way DWD; Lyu et al. 2017), a tensor logistic regression based on the Tucker decomposition (Tucker; Li et al. 2018), a regularized matrix regression (RMR; Zhou & Li 2014), and the covariate-adjusted tensor classification in high-dimensions (CATCH; Pan et al. 2019). We focus on the classification error rates of these methods. Therefore, Bayes' error is also reported.

5.2 Simulations under the DATE model

Unless otherwise specified, we generate data from the DATE model as follows:

$$(\mathbf{X} \mid Y = k) \sim \boldsymbol{\mu}_k + \llbracket \mathbf{Z}; \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_M \rrbracket, \sum_{k=1}^K (n_k/n) \boldsymbol{\mu}_k = \mathbf{0}, \quad (5.9)$$

where \mathbf{Z} consists of independent $N(0, 1)$ random variables, such that $\mathbf{X} \mid (Y = k) \sim TN(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_M)$. We first let $\boldsymbol{\mu}_k^* = \llbracket \boldsymbol{\Theta}_k; \boldsymbol{\Gamma}_1, \dots, \boldsymbol{\Gamma}_M \rrbracket$ for some randomly generated $\boldsymbol{\Theta}_k \in \mathbb{R}^{u_1 \times \dots \times u_M}$ with Uniform(0, 1) elements, and let $\bar{\boldsymbol{\mu}}^* = \sum_{k=1}^K (n_k/n) \boldsymbol{\mu}_k^*$. Then, the mean parameter $\boldsymbol{\mu}_k = \boldsymbol{\mu}_k^* - \bar{\boldsymbol{\mu}}^*$. Let $\boldsymbol{\Sigma}_m^* = \boldsymbol{\Gamma}_m \boldsymbol{\Omega}_m \boldsymbol{\Gamma}_m^T + \boldsymbol{\Gamma}_{0m} \boldsymbol{\Omega}_{0m} \boldsymbol{\Gamma}_{0m}^T$. The covariance matrix $\boldsymbol{\Sigma}_m = \sigma^2 \times \boldsymbol{\Sigma}_m^* / \|\boldsymbol{\Sigma}_m^*\|_F$,

where the scalar $\sigma^2 > 0$ is chosen differently for each model to control Bayes' classification error in a reasonable range. The envelope basis matrices Γ_m , for $m = 1, \dots, M$, are generated with Uniform(0, 1) elements and then orthogonalized. Three types of covariance are considered:

- (C1) $\Omega_m = \mathbf{I}_{u_m}$ and $\Omega_{0m} = 0.01\mathbf{I}_{p_m - u_m}$;
- (C2) $\Omega_m = 0.1\mathbf{I}_{u_j}$ and $\Omega_{0m} = \mathbf{I}_{p_m - u_m}$;
- (C3) $\Omega_m = \mathbf{O}_m \mathbf{D}_m \mathbf{O}_m^T$ and $\Omega_{0m} = \mathbf{O}_{0m} \mathbf{D}_{0m} \mathbf{O}_{0m}^T$, where $\mathbf{O}_m \in \mathbb{R}^{u_m \times u_m}$ and $\mathbf{O}_{0m} \in \mathbb{R}^{(p_m - u_m) \times (p_m - u_m)}$ are randomly generated orthogonal matrices, \mathbf{D}_{0m} is a diagonal matrix with elements $\exp(k_{m,1}, \dots, k_{m,p_m - u_m})$, where $k_{m,1}, \dots, k_{m,p_m - u_m}$ are $p_m - u_m$ evenly spaced numbers between -10 and m , and \mathbf{D}_m is a diagonal matrix to be specified later.

The following four DATE models are considered. For each training set with sample size n , we evaluate the predictive performance of each method on a testing set with sample size of $10n$. The results are averaged over 100 replications. For the tuning parameters in each model, we use the true envelope dimension $\{u_m\}_{m=1}^M$ for our method and for the Tucker rank. The m -way DWD uses rank $\prod_{m=1}^M u_m$. All tuning parameters in the penalized methods (CATCH, ℓ_1 -GLM, ℓ_1 -FDA, and RMR) are chosen using five-fold cross-validation.

- (M1) Matrix predictor with binary response, $M = K = 2$. We generate training data with $n = 200$ observations. Let $p_1 = 80$, $p_2 = 20$, $u_1 = 4$, $u_2 = 2$, and $n_1 = n_2 = n/2$. The parameter σ^2 is 1, 40, and 3 for

covariance (C1)–(C3), respectively. For (C3), \mathbf{D}_m is a diagonal matrix with u_m elements $(5, 5^2, \dots)$.

- (M2) The true parameter $\mathbf{B} \in \mathbb{R}^{p_1 \times p_2}$ is constructed such that we can visualize the estimates directly (see Figure S1 in Supplementary Material).

Let $n = 300$, $p_1 = p_2 = 64$, $u_1 = u_2 = 2$, and $n_1 = n_2 = n/2$. The parameter σ^2 is 0.1, 5, and 0.13 for covariances (C1)–(C3), respectively.

For (C3), \mathbf{D}_m is a diagonal matrix with u_m elements (e, e^2, \dots) .

- (M3) Similar to (M1), but with $K = 4$. Let $n = 300$, $u_1 = u_2 = 5$, $p_1 = p_2 = 50$, and $n_1 = n_2 = n_3 = n/4$. The parameter σ^2 is 1.5, 40, and 2.5 for covariances (C1)–(C3), respectively.

- (M4) Similar to (M1), but with $M = K = 3$. Let $p_1 \times p_2 \times p_3 = 20 \times 30 \times 40$, $u = (2, 3, 4)$, and $n_1 = 90$, $n_2 = 60$, $n_3 = 150$. The parameter σ^2 is 1.3, 40, and 2 for covariances (C1)–(C3), respectively.

The results for the above four models and the three covariance structures are summarized in Table 1. Note that some binary and matrix classification methods, such as m -way DWD, Tucker, and RMR, cannot be applied to three-way tensor data in M3 and M4. In addition, Tucker is not applicable for M1, because the sample size is too small.

Under the covariance structure (C1), the material variation in the predictor is much larger than the immaterial variation. The setting is thus not challenging, and most of the methods work well. Under (C2), the immaterial variation dominates, and most methods fail to identify the weak signals. The only exception is

5. NUMERICAL STUDIES²⁷

Model	M1			M2		
	(C1)	(C2)	(C3)	(C1)	(C2)	(C3)
Bayes	11.53	14.03	14.63	13.74	16.29	14.08
DATE-L	13.25 (0.15)	16.74 (0.10)	16.36 (0.15)	14.07 (0.10)	17.01 (0.19)	14.26 (0.06)
DATE-D	12.45(0.09)	50.05 (0.10)	16.12 (0.13)	14.06 (0.08)	49.91 (0.08)	15.14 (0.07)
TDA	33.34 (0.13)	36.08 (0.12)	35.95 (0.14)	38.13 (0.11)	40.23 (0.11)	38.46 (0.11)
m -way DWD	12.36 (0.09)	50.69 (0.11)	21.73 (0.17)	13.99 (0.07)	49.98 (0.09)	14.66 (0.07)
CATCH	13.53 (0.13)	49.91 (0.11)	19.13 (0.23)	14.58 (0.08)	49.93 (0.10)	15.37 (0.10)
Tucker	-	-	-	40.49 (0.22)	48.43 (0.15)	42.23 (0.22)
RMR	12.05 (0.09)	49.35 (0.11)	18.17 (0.13)	14.02 (0.07)	49.86 (0.10)	14.72 (0.08)
DLDA	12.74 (0.09)	49.99 (0.12)	26.52 (0.22)	14.91 (0.08)	49.74 (0.09)	19.90 (0.13)
ℓ_1 -GLM	13.65 (0.16)	49.90 (0.06)	17.88 (0.12)	15.53 (0.09)	50.00 (0.06)	16.79 (0.10)
ℓ_1 -FDA	12.74 (0.09)	49.95 (0.09)	26.52 (0.22)	14.91 (0.08)	49.87 (0.06)	19.90 (0.13)
Model	M3			M4		
	(C1)	(C2)	(C3)	(C1)	(C2)	(C3)
Bayes	16.45	12.04	16.92	11.54	12.69	11.28
DATE-L	20.01 (0.18)	14.47 (0.07)	19.75 (0.13)	19.75 (0.43)	22.05 (0.25)	15.63 (0.27)
DATE-D	19.72(0.11)	74.67 (0.16)	21.79 (0.10)	14.85 (0.09)	56.72 (0.10)	15.14 (0.13)
TDA	54.35 (0.10)	49.21 (0.11)	55.09 (0.11)	49.88 (0.01)	49.90 (0.01)	49.85 (0.01)
CATCH	21.59 (0.12)	74.88 (0.07)	25.90 (0.19)	15.16 (0.10)	53.91 (0.25)	25.14 (0.27)
DLDA	31.00 (0.13)	75.16 (0.09)	41.77 (0.15)	13.47 (0.07)	50.19 (0.02)	20.56 (0.11)
ℓ_1 -GLM	21.43 (0.10)	75.00 (0.05)	25.64 (0.13)	15.55 (0.09)	50.34 (0.48)	22.90 (0.13)
ℓ_1 -FDA	18.65 (0.08)	74.83 (0.07)	27.74 (0.12)	13.47 (0.07)	50.02 (0.01)	20.56 (0.11)

Table 1: Averaged classification error (%) and standard error (in parentheses), calculated over 100 replicates.

DATE-L, which effectively identifies that Ω_m contains the small eigenvalues in Σ_m . Finally, covariance structure (C3) is between the two extremes of (C1) and (C2). Its complex covariance structure favors both DATE-L and DATE-D over other methods.

From Table 1, DATE-L is either the best or very close to the best for all the models considered. Moreover, it is the only method that works well under the

covariance structure (C2). Although DATE-D is not a likelihood-based method, it has very good finite-sample performance that is similar to that of DATE-L under the (C1) and (C3) covariance structures. This is an encouraging result for DATE-D, because it is a much faster and simpler estimation method for the tensor envelopes. For the more complex covariance (C3), DATE-L and DATE-D outperform the other methods, improving estimation significantly.

Comparing models M3 (multi-class response and matrix predictor) and M4 (multi-class response and tensor predictor) with model M1 (binary response and matrix predictor), the advantages of DATE over TDA (and other methods) is more significant when $K, M > 2$. In model M2, several estimators (DATE-L, Tucker, and RMR) of $\mathbf{B} \in \mathbb{R}^{64 \times 64}$ are visualized in Figure S1 in the Supplementary Material, showing that DATE-L clearly provides much better parameter estimates than its competitors do under all three covariance structures.

5.3 Violation of DATE model assumptions

We consider the following models in which the DATE assumptions are violated.

- (Heavy-tail distribution) We first consider a model in which the data are generated from a multivariate t -distribution. The model is the same as (M1), except that we set $p_1 = p_2 = 20$, $u_1 = u_2 = 2$, and each element in \mathbf{Z} is generated independently from a Student's t -distribution with degree of freedom 4. The parameter σ^2 is 0.4, 40, and 3 for covariances (C1)–(C3), respectively. For (C3), \mathbf{D}_m is diagonal with u_m elements $(5, 5^2, \dots)$.

5. NUMERICAL STUDIES₂₉

- (TDA models) We consider two TDA models from (3.1), where no envelope assumptions are imposed on $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_m$. The true envelope dimension is $u_m = p_m$. We set $p_1 = p_2 = 20$, $K = 2$, and $n_1 = n_2 = 100$:
 - (TDA1) Let $\boldsymbol{\mu}_k = \llbracket \boldsymbol{\Theta}_k; \boldsymbol{\Gamma}_1, \boldsymbol{\Gamma}_2 \rrbracket$ for a randomly generated $\boldsymbol{\Theta}_k \in \mathbb{R}^{2 \times 2}$ with Uniform(0, 1) elements and randomly generated basis matrices $\boldsymbol{\Gamma}_m \in \mathbb{R}^{20 \times 2}$, \mathbf{D}_m be a diagonal matrix with elements evenly spaced between 0.3 and 3, and $\mathbf{O}_m \in \mathbb{R}^{p_m \times p_m}$ be a randomly generated orthogonal matrix. Then, we let $\boldsymbol{\Sigma}_m = \sigma^2 \boldsymbol{\Sigma}_m^* / \|\boldsymbol{\Sigma}_m^*\|_F$, where $\boldsymbol{\Sigma}_m^* = \mathbf{O}_m \mathbf{D}_m \mathbf{O}_m^T$ and $\sigma^2 = 1.2$.
 - (TDA2) Each element of $\boldsymbol{\mu}_k$ is generated randomly from Uniform(0.2, 1) and $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2 = 2.5\text{AR}(0.3)$, where $\text{AR}(\rho)$ represents a covariance matrix with the (i, j) -th element to be $\rho^{|i-j|}$.

We also construct three simulation models from our competing methods m -way DWD (Lyu et al. 2017), Tucker logistic regression (Li et al. 2018), and CATCH (Pan et al. 2019). This allows us to better understand how the two DATE methods perform under model misspecification.

- (DWD model) Let $p_1 \times p_2 = 20 \times 10$ and $n_1 = n_2 = 50$. For each training data set, the vectorized samples are generated from a multivariate normal distribution $N(\text{vec}(\boldsymbol{\mu}_k), \boldsymbol{\Sigma}_{ek})$, with $\boldsymbol{\Sigma}_{ek} = \sigma_{ek}^2 \mathbf{I}$. Let $\boldsymbol{\mu}_1 = \mathbf{0}$ and $\boldsymbol{\mu}_2 = \mathbf{v} \otimes \mathbf{w}$, which corresponds to a rank-1 DWD model, where \mathbf{v} and \mathbf{w} are generated from multivariate normal distributions with mean zero and

variances $\sigma_v^2 \mathbf{I}$ and $\sigma_w^2 \mathbf{I}$, respectively. We set $\sigma_{e_1}^2 = 1.75$, $\sigma_{e_2}^2 = 2$, and $\sigma_v^2 = \sigma_w^2 = 0.2$.

- (Tucker model) For the Tucker logistic regression model, the regression coefficient $\mathbf{B} \in \mathbb{R}^{64 \times 64}$ is the same as in M2, the predictors $\mathbf{X}_i \in \mathbb{R}^{64 \times 64}$, for $i = 1, \dots, n$, are randomly generated with all elements being independent standard normal. The binary response Y is generated from a binomial distribution with probability $\{1 + \exp(-\langle \mathbf{B}, \mathbf{X} \rangle)\}^{-1}$. The training sample size $n = 500$ and the testing sample size is 1000.
- (CATCH model) TDA model (3.1) with $p_1 = p_2 = p_3 = 20$, $n_1 = n_2 = 100$, and $K = 2$. The parameter \mathbf{B} is sparse with nonzeros at the $3 \times 3 \times 3$ sub-tensor generated from $\text{Uniform}(0, 0.5)$ independently. Let \mathbf{D} be a diagonal matrix with diagonal elements evenly between 3 and 0.3. We set $\Sigma_1 = \mathbf{D}$, $\Sigma_2 = (\mathbf{D} + \text{AR}(0.3))/2$, and $\Sigma_3 = (\mathbf{D} + \text{CS}(0.3))/2$, where $\text{CS}(\rho)$ is a matrix with diagonals that are ones, and off-diagonals that are equal to ρ 's.

For the results of DATE-D and DATE-L presented in Table 2, we use $u_1 = u_2 = 2$ for the t -distribution model, and cross-validation to select the dimensions for all the other models. DATE-L is the best or very close to the best method, overall, for all the models in this section, and DATE-D is fairly competitive for most of the models. For the t -distribution model, DATE-L is among the best methods for all the covariance structures, and the only method that works well for covariance (C2). DATE-D performs similarly to DATE-L for covari-

ances (C1) and (C3). Thus, the two proposed DATE methods are not sensitive to nonnormal heavy-tailed distributions. In the TDA1 model, the parameterization (3.5) for the covariance matrices is violated, and the mean parameter still has a low-rank structure. Note that DATE-L can still find a low-dimensional subspace such that the projected data are informative, and provides better classification results than those of TDA. We visualize the classification errors of the various methods in Figure 1, where we vary the envelope dimension. It is clear that DATE-L has superior classification accuracy when the input dimension is between 5 and 15. From Figure 1, as the input envelope dimensions increase, DATE-D improves and reaches the same results as TDA when $u_m = p_m$. In the TDA2 model, the parameterizations (3.5) and (3.6) are both violated. In this case, cross-validation returns a dimension equal to p_m or close to p_m for both DATE-L and DATE-D, resulting in almost identical performance to that of TDA.

Not surprisingly, the m -way DWD, Tucker logistic regression, and CATCH perform best under their own respective models. It is encouraging to note that DATE-L exhibits competitive performance with the best methods, and is better overall, demonstrating the flexibility and effectiveness of our DATE-L estimator. DATE-D performs well for the CATCH model and the DWD model, but, in general, is less effective than DATE-L. We believe that DATE-L, when applicable, is probably more robust than DATE-D under model misspecification.

5. NUMERICAL STUDIES³²

Model	t distribution			TDA		DWD	Tucker	CATCH
	(C1)	(C2)	(C3)	(TDA1)	(TDA2)			
Bayes	16.10	16.10	12.22	6.50	6.18	11.82	-	5.63
DATE-L	16.66 (0.14)	16.81 (0.13)	12.73 (0.07)	8.92 (0.11)	13.75 (0.15)	23.82 (0.24)	30.99 (0.22)	11.71 (0.27)
DATE-D	16.64 (0.09)	50.14 (0.08)	12.84 (0.08)	11.16 (0.13)	13.46 (0.14)	30.93 (0.20)	46.14 (0.16)	10.54 (0.12)
TDA	29.34 (0.14)	29.33 (0.14)	24.59(0.13)	10.58 (0.09)	13.13 (0.09)	30.67 (0.20)	41.51 (0.15)	35.16 (0.12)
<i>m</i> -way DWD	16.70 (0.09)	50.06 (0.10)	13.97(0.08)	10.49 (0.09)	36.88 (0.23)	21.16 (0.18)	48.34 (0.55)	-
CATCH	17.34 (0.12)	49.49 (0.10)	14.14 (0.11)	14.06 (0.13)	22.81 (0.19)	33.42 (0.41)	42.76 (0.19)	7.44 (0.09)
Tucker	38.48 (0.22)	40.36 (0.30)	34.09 (0.27)	27.91 (0.28)	34.84 (0.27)	41.84 (0.32)	28.18 (0.25)	-
RMR	16.40 (0.08)	49.34 (0.12)	14.09 (0.09)	10.44 (0.08)	21.75 (0.14)	33.52 (0.20)	39.63 (0.16)	-
DLDA	16.65 (0.08)	49.49 (0.12)	16.54 (0.15)	12.95 (0.09)	22.81 (0.14)	39.28 (0.23)	41.37 (0.15)	34.72 (0.14)
ℓ_1 -GLM	17.81 (0.08)	49.93 (0.07)	13.97 (0.15)	20.21 (0.17)	30.14 (0.20)	38.42 (0.54)	46.70 (0.29)	11.81 (0.14)
ℓ_1 -FDA	16.65 (0.10)	49.84 (0.07)	16.54 (0.10)	12.95 (0.09)	22.81 (0.14)	29.23 (0.21)	41.37 (0.15)	34.71 (0.14)

Table 2: The averaged error rates and associated standard errors over 100 replicates.

5.4 Data sets with higher dimensions

We consider simulations in which the tensor dimension $p = \prod_m p_m$ is much larger than the sample size n . We vary the sample size of the training set from 50 to 400, and set the sample size of the testing set to 2000. The settings for this model are analogous to those for M1, but $p_1 = p_2 = 200$ and $u_1 = u_2 = 5$. We use covariance (C3) with $\mathbf{D}_m = \mathbf{I}_{u_m}$ and $\sigma^2 = 3$.

From Table 3, DATE-D still performs well even when $p_1 \times p_2$ is much larger than n (e.g., $n=100$). Note that DATE-D outperforms DATE-L, especially when n is small. DATE-L is not accurate and is less stable when the dimensions of the predictors are much larger than n , partially because of the nonconvex optimization. This simulation model provides encouraging evidence that DATE-D can

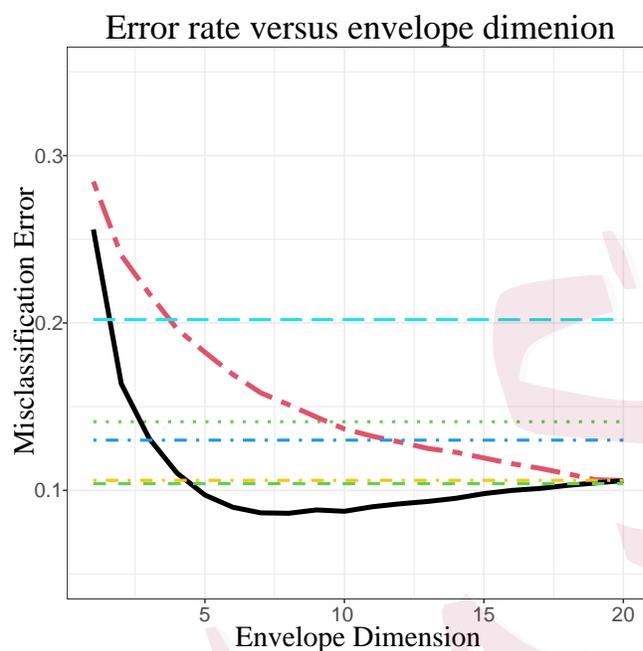


Figure 1: Model (TDA1): Classification error versus the input envelope dimensions $u_1 = u_2$. The black curve shows DATE-L and the red dashed curve shows DATE-D. From top to bottom, the horizontal dashed lines show l_1 -GLM, CATCH, DLDA, TDA, and RMR, respectively.

be applied in high dimensions. The results support those of Theorem 4, which states that DATE-D can be consistent, even when p_m goes to infinity faster than n does. In our experience, DATE-D can handle cases with much larger dimensions and is computationally efficient, because it involves only matrix multiplications and eigen-decompositions.

5. NUMERICAL STUDIES³⁴

n	$\{p_m\} = \{200, 200\}, n_1 = n_2 = n/2$			
	50	100	200	400
Bayes	7.89	7.89	7.89	7.89
DATE-D	24.33 (0.42)	15.85 (0.23)	11.63 (0.13)	9.74 (0.08)
DATE-L	46.70 (0.56)	30.85 (0.88)	15.38(0.30)	10.65 (0.11)
TDA	47.12 (0.12)	45.93 (0.13)	44.50 (0.11)	41.92 (0.11)
CATCH	39.17 (0.42)	29.30 (0.32)	20.75 (0.24)	14.91 (0.13)
ℓ_1 -GLM	43.10 (0.52)	31.54 (0.52)	22.78 (0.16)	17.59 (0.11)
ℓ_1 -FDA	22.63 (0.19)	17.15 (0.13)	13.18 (0.10)	10.62 (0.08)

Table 3: The averaged error rates and associated standard errors over 100 replicates.

5.5 Real-data examples

The first data set is from a colorimetric sensor array (CSA) study, where chemical dyes are used to transform smell into optical composite signals (Zhong & Suslick 2015). The experiments used a colorimetric sensor array separately at Immediately Dangerous to Life or Health (IDLH) and Permissible Exposure Level (PEL) concentrations of the $K = 21$ chemical toxicants. The dimension of the predictor is 36×3 , and the total sample size is $n = 7 \times K = 147$. For each of the data sets, IDLH and PEL, we perform 100 repeated training/testing splits, 126 as training and 21 as testing, because each class has only seven samples. The tuning parameters for all methods are based on cross-validation. For DATE-L, we select the dimensions $u_1 = 8$ and $u_2 = 2$ for the IDLH data set, and $u_1 = 7$ and $u_2 = 3$ for the PEL data set. For DATE-D, we select $u_1 = 7$ and $u_2 = 2$ for the IDLH data set, and $u_1 = 9$ and $u_2 = 3$ for the PEL data set. The results are summarized in Table 4. Because of a large number of classes

5. NUMERICAL STUDIES³⁵

and very low sample sizes per class, many methods are not applicable. It is clear that the DATE-L and DATE-D methods achieve better classification results than the other methods do on this data set. In particular, DATE-L, DATE-D, CATCH, and ℓ_1 -FDA achieve perfect classification in the IDLH setting. In the PEL data set, the classification becomes more difficult. Here, DATE-L achieves the best classification, followed by DATE-D.

	CSA-IDLH	CSA-PEL	GT (LOO)	GT (10-fold CV)
DATE-L	0 (0)	1.24 (0.28)	9.43	12.20 (0.30)
DATE-D	0 (0)	2.24 (0.26)	11.32	14.77 (0.26)
TDA	-	-	-	-
<i>m</i> -way DWD	-	-	16.98	17.63 (0.33)
CATCH	0 (0)	4.03 (0.43)	16.98	20.50 (0.31)
RMR	-	-	15.09	20.07 (0.23)
LDA	4.81 (0.52)	17.81 (0.78)	32.08	30.23 (0.27)
ℓ_1 -GLM	0.57 (0.17)	16.89 (0.50)	26.42	26.88 (0.43)
ℓ_1 -FDA	0 (0)	6.62(0.55)	32.08	30.55 (0.27)

Table 4: The classification errors, averaged over different training-testing sample splits: seven-fold cross-validation for the CSA data, and leave-one-out (LOO) and 10-fold cross-validation for the GT data.

The second study is the Gene Time (GT) study of Baranzini et al. (2004), who collected gene expressions from patients suffering from multiple sclerosis (MS). Fifty-three patients treated with recombinant human interferon beta ($rIFN\beta$) are followed at six time points with 76-gene expression data, resulting in tensor data of dimension $p_1 \times p_2 = 76 \times 7$ and $n = 53$. This is a binary classification problem. The two classes are patients who respond well and those who respond poorly to interferon beta. Based on cross-validation, we select $u_1 = 5$

6. CONCLUSION

and $u_2 = 1$ for the DATE-L method, and $u_1 = 7$ and $u_2 = 1$ for the DATE-D method. We compare the classification errors for leave-one-out (LOO) and 10-fold cross-validation. For the 10-fold cross-validation, we repeat 100 times, and report the average classification errors and standard errors. The results are summarized in Table 4. Again, DATE-L has the lowest error rate, followed by DATE-D, based on both leave-one-out and 10-fold cross-validation. The improvement of the DATE methods over the other methods is quite substantial. These encouraging results indicate that the DATE methods can capture information in both the parameter B_k and the covariance matrices.

6. Conclusion

We have developed a parsimonious tensor discriminant analysis model based on tensor envelopes. A likelihood-based estimator is derived from the tensor normal likelihood, and is shown to be effective in practice and robust to model assumption violations. When the tensor dimension is very high and the likelihood-based estimator becomes infeasible, a fast decomposition-based estimator can be applied with theoretical guarantees.

The estimators can be extended to the covariate-adjusted tensor classification framework of Pan et al. (2019). Details of this extension, including the formulation of the DATE estimators and the derivations of the MLEs, as well as simulations and a real-data example are included in the Supplementary Material, Section S1.

7. SUPPLEMENTARY MATERIALS

7. Supplementary Materials

The online supplementary material contains additional numerical results, implementation details, and proofs.

References

- Baranzini, S. E., Mousavi, P., Rio, J., Caillier, S. J., Stillman, A., Villoslada, P., Wyatt, M. M., Comabella, M., Greller, L. D. & Somogyi, R. (2004), 'Transcription-based prediction of response to $\text{ifn}\beta$ using supervised computational methods', *PLoS Biol* **3**(1), e2.
- Bi, Y. & Jeske, D. R. (2010), 'The efficiency of logistic regression compared to normal discriminant analysis under class-conditional classification noise', *Journal of Multivariate Analysis* **101**(7), 1622–1637.
- Chi, E. C. & Kolda, T. G. (2012), 'On tensors, sparsity, and nonnegative factorizations', *SIAM Journal on Matrix Analysis and Applications* **33**(4), 1272–1299.
- Cichocki, A., Mandic, D., De Lathauwer, L., Zhou, G., Zhao, Q., Caiafa, C. & Phan, H. A. (2015), 'Tensor decompositions for signal processing applications: From two-way to multiway component analysis', *IEEE signal processing magazine* **32**(2), 145–163.
- Conway, J. B. (2013), *A course in functional analysis*, Vol. 96, Springer Science & Business Media.
- Cook, R. D. (2018), *An Introduction to Envelopes: Dimension Reduction for Efficient Estimation in Multivariate Statistics*, John Wiley & Sons.
- Cook, R. D. & Zhang, X. (2015), 'Foundations for envelope models and methods', *Journal of the American Statistical Association* **110**(510), 599–611.
- Cook, R. D. & Zhang, X. (2016), 'Algorithms for envelope estimation', *Journal of Computational and Graphical Statistics* **25**(1), 284–300.
- Cook, R., Li, B. & Chiaromonte, F. (2010), 'Envelope models for parsimonious and efficient multivariate

REFERENCES

- linear regression', *Statistica Sinica* **20**(3), 927–960.
- Dudoit, S., Fridlyand, J. & Speed, T. P. (2002), 'Comparison of discrimination methods for the classification of tumors using gene expression data', *Journal of the American Statistical Association* **97**(457), 77–87.
- Efron, B. (1975), 'The efficiency of logistic regression compared to normal discriminant analysis', *Journal of the American Statistical Association* **70**(352), 892–898.
- Fisher, R. A. (1936), 'The use of multiple measurements in taxonomic problems', *Annals of Eugenics* **7**(2), 179–188.
- Friedman, J., Hastie, T. & Tibshirani, R. (2010), 'Regularization paths for generalized linear models via coordinate descent', *Journal of Statistical Software* **33**(1), 1–22.
- Gahrooei, M. R., Yan, H., Paynabar, K. & Shi, J. (2020), 'Multiple tensor-on-tensor regression: An approach for modeling processes with heterogeneous sources of data', *Technometrics* **63**(2), 1–23.
- Gupta, A. K. & Nagar, D. K. (2018), *Matrix variate distributions*, Chapman and Hall/CRC.
- Hitchcock, F. L. (1927), 'The expression of a tensor or a polyadic as a sum of products', *Journal of Mathematics and Physics* **6**(1-4), 164–189.
- Hoff, P. D. (2015), 'Multilinear tensor regression for longitudinal relational data', *The Annals of Applied Statistics* **9**(3), 1169.
- Hore, V., Viñuela, A., Buil, A., Knight, J., McCarthy, M. I., Small, K. & Marchini, J. (2016), 'Tensor decomposition for multiple-tissue gene expression experiments', *Nature Genetics* **48**(9), 1094.
- Kolda, T. G. & Bader, B. W. (2009), 'Tensor decompositions and applications', *SIAM Review* **51**(3), 455–500.
- Li, L. & Zhang, X. (2017), 'Parsimonious tensor response regression', *Journal of the American Statistical Association* **112**(519), 1131–1146.

REFERENCES

- Li, P. & Maiti, T. (2019), Universal consistency of support tensor machine, in '2019 IEEE International Conference on Data Science and Advanced Analytics (DSAA)', IEEE, pp. 608–609.
- Li, Q. & Schonfeld, D. (2014), 'Multilinear discriminant analysis for higher-order tensor data classification', *IEEE transactions on pattern analysis and machine intelligence* **36**(12), 2524–2537.
- Li, X., Xu, D., Zhou, H. & Li, L. (2018), 'Tucker tensor regression and neuroimaging analysis', *Statistics in Biosciences* **10**(3), 520–545.
- Lock, E. F. (2018), 'Tensor-on-tensor regression', *Journal of Computational and Graphical Statistics* **27**(3), 638–647.
- Lyu, T., Lock, E. F. & Eberly, L. E. (2017), 'Discriminating sample groups with multi-way data', *Biostatistics* **18**(3), 434–450.
- Manceur, A. M. & Dutilleul, P. (2013), 'Maximum likelihood estimation for the tensor normal distribution: Algorithm, minimum sample size, and empirical bias and dispersion', *Journal of Computational and Applied Mathematics* **239**, 37–49.
- Molstad, A. J. & Rothman, A. J. (2019), 'A penalized likelihood method for classification with matrix-valued predictors', *Journal of Computational and Graphical Statistics* **28**(1), 11–22.
- Pan, Y., Mai, Q. & Zhang, X. (2019), 'Covariate-adjusted tensor classification in high dimensions', *Journal of the American Statistical Association* **114**(527), 1305–1319.
- Raskutti, G., Yuan, M., Chen, H. et al. (2019), 'Convex regularization for high-dimensional multiresponse tensor regression', *The Annals of Statistics* **47**(3), 1554–1584.
- Sun, W. W. & Li, L. (2017), 'Store: sparse tensor response regression and neuroimaging analysis', *The Journal of Machine Learning Research* **18**(1), 4908–4944.
- Tucker, L. R. (1966), 'Some mathematical notes on three-mode factor analysis', *Psychometrika* **31**(3), 279–311.

REFERENCES

- Wang, X., Zhu, H. & Initiative, A. D. N. (2017), ‘Generalized scalar-on-image regression models via total variation’, *Journal of the American Statistical Association* **112**(519), 1156–1168.
- Wang, Y., Meng, D. & Yuan, M. (2018), ‘Sparse recovery: from vectors to tensors’, *National Science Review* **5**(5), 756–767.
- Witten, D. M. & Tibshirani, R. (2011), ‘Penalized classification using fisher’s linear discriminant’, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **73**(5), 753–772.
- Yan, H., Paynabar, K. & Pacella, M. (2019), ‘Structured point cloud data analysis via regularized tensor regression for process modeling and optimization’, *Technometrics* **61**(3), 385–395.
- Yan, S., Xu, D., Yang, Q., Zhang, L., Tang, X. & Zhang, H.-J. (2006), ‘Multilinear discriminant analysis for face recognition’, *IEEE Transactions on Image Processing* **16**(1), 212–220.
- Ye, J., Janardan, R. & Li, Q. (2004), ‘Two-dimensional linear discriminant analysis’, *Advances in neural information processing systems* **17**, 1569–1576.
- Zhang, A. (2019), ‘Cross: Efficient low-rank tensor completion’, *The Annals of Statistics* **47**(2), 936–964.
- Zhang, A. & Xia, D. (2018), ‘Tensor svd: Statistical and computational limits’, *IEEE Transactions on Information Theory* **64**(11), 7311–7338.
- Zhang, X., Deng, K. & Mai, Q. (2021), ‘Envelopes and principal component regression’, *Manuscript* .
- Zhang, X. & Li, L. (2017), ‘Tensor envelope partial least-squares regression’, *Technometrics* **59**(4), 426–436.
- Zhang, X. & Mai, Q. (2018), ‘Model-free envelope dimension selection’, *Electronic Journal of Statistics* **12**(2), 2193–2216.
- Zhang, X. & Mai, Q. (2019), ‘Efficient integration of sufficient dimension reduction and prediction in discriminant analysis’, *Technometrics* **61**(2), 259–272.
- Zhong, W. & Suslick, K. S. (2015), ‘Matrix discriminant analysis with application to colorimetric sensor

REFERENCES

array data', *Technometrics* **57**(4), 524–534.

Zhou, H. & Li, L. (2014), 'Regularized matrix regression', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **76**(2), 463–483.

Zhou, H., Li, L. & Zhu, H. (2013), 'Tensor regression with applications in neuroimaging data analysis', *Journal of the American Statistical Association* **108**(502), 540–552.