

UCLA

UCLA Previously Published Works

Title

Tensor canonical correlation analysis

Permalink

<https://escholarship.org/uc/item/565058m6>

Journal

Stat, 8(1)

ISSN

0038-9986

Authors

Min, Eun Jeong

Chi, Eric C

Zhou, Hua

Publication Date

2019

DOI

10.1002/sta4.253

Peer reviewed



Published in final edited form as:

Stat. 2020 ; 8(1): . doi:10.1002/sta4.253.

Tensor canonical correlation analysis

Eun Jeong Min¹, Eric C. Chi², Hua Zhou³

¹Department of Biostatistics, Epidemiology and Informatics, University of Pennsylvania, Philadelphia, 19104, PA, U.S.A.

²Department of Statistics, North Carolina State University, Raleigh, 27695, NC, U.S.A.

³Department of Biostatistics, University of California, Los Angeles, Los Angeles, 90095, CA, U.S.A.

Abstract

Canonical correlation analysis (CCA) is a multivariate analysis technique for estimating a linear relationship between two sets of measurements. Modern acquisition technologies, for example, those arising in neuroimaging and remote sensing, produce data in the form of multidimensional arrays or tensors. Classic CCA is not appropriate for dealing with tensor data due to the multidimensional structure and ultrahigh dimensionality of such modern data. In this paper, we present tensor CCA (TCCA) to discover relationships between two tensors while simultaneously preserving multidimensional structure of the tensors and utilizing substantially fewer parameters. Furthermore, we show how to employ a parsimonious covariance structure to gain additional stability and efficiency. We delineate population and sample problems for each model and propose efficient estimation algorithms with global convergence guarantees. Also we describe a probabilistic model for TCCA that enables the generation of synthetic data with desired canonical variates and correlations. Simulation studies illustrate the performance of our methods.

Keywords

block coordinate ascent; CP decomposition; multidimensional array data

1 | INTRODUCTION

Canonical correlation analysis (CCA) is a classic statistical method for identifying associations between two sets of measurements (Hotelling, 1936). Specifically, CCA identifies a pair of coefficient vectors, one for each set of measurements, such that the correlation between the corresponding linear combinations of variables from each set is maximized. By default, CCA applies when each observation consists of a pair of vector covariates. In many modern data analysis problems, however, each observation may consist more generally of a pair of multidimensional arrays or tensors. For example, in imaging

Correspondence: Eric C. Chi, Department of Statistics, North Carolina State University, Raleigh, NC 27695, U.S.A. eric_chi@ncsu.edu.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of the article.

genetics, to identify genetic variants that can best capture and explain phenotypic variations in brain function and structure, Stein et al. (2010) studied $p = 448,293$ single-nucleotide polymorphisms and $q = 31,622$ voxels in brain images, on $n = 740$ individuals. A naive approach to dealing with tensor-valued data would be to reshape tensor covariates into vectors and then apply standard CCA. There are, however, two serious drawbacks to doing so. First, structural information in tensors is discarded through vectorization. Second, the resulting vectors consist of a prohibitively large number of parameters. In the imaging genetics problem (Stein et al., 2010), vectorizing the voxel intensity measurements disregards the spatial correlation among neighbouring voxels. Moreover, applying standard CCA to vectors of single-nucleotide polymorphism covariates and vectorized brain images would require estimating nearly half a million parameters using fewer than a thousand observations.

In light of these issues, there have been extensions of CCA to handle special cases of tensor-valued data (Lee & Choi, 2007; Wang, 2010; Yan, Zheng, Zhou, & Zhao, 2012; Gang, Yong, Yan-Lei, & Jing, 2011; Lu, 2013; Wang, Yan, Sun, Zhao, & Fu, 2016). Although they have exhibited good empirical performance in some applications, there remains no clear population models underlying these sample based heuristics. To address this gap in the literature, we introduce a novel statistical model for tensor CCA (TCCA). We summarize at a high level our formulation and its contributions:

- We propose a TCCA population model that imposes the CANDECOMP/PARAFAC (CP) decomposition (Carroll & Chang, 1970; Harshman, 1970) structure on canonical tensors. This population model enforces model parsimony and enables efficient estimation.
- We propose a refinement of TCCA, which assumes a separable covariance structure (scTCCA). This refinement enables efficient estimation of large covariance matrices of tensor-valued data.
- We derive convenient representations of the covariance between linear combinations of two random tensors under the unstructured and separable covariance structured assumptions.
- We develop efficient estimation algorithms for TCCA and scTCCA, both based on block coordinate ascent, which leverage these efficient representations. Each step of both algorithms solves a substantially lower dimensional CCA problem; thus, both algorithms can be easily implemented using any standard solvers for the CCA problem. Moreover, we prove global convergence guarantees of both estimation algorithms under modest regularity conditions.
- We develop simple modifications to the TCCA and scTCCA estimation algorithms to incorporate recovery of sparse canonical correlation tensors to improve interpretability of the estimated models.
- Finally, we extend the probabilistic interpretation of CCA by Bach and Jordan (2006) to TCCA. This extension leads to a probabilistic model for generating datasets with specified canonical correlation and variates.

The remainder of the paper is organized as follows. In Section 2, we review tensor notation and basic operations used in this paper. In Sections 3 and 4, we propose our two TCCA methods: TCCA and scTCCA. In Section 5, we describe a modification of the estimation algorithms for TCCA and scTCCA for sparse models. In Section 6, we introduce the probabilistic TCCA model. In Section 7, we describe the numerical experiment results. In Section 8, we conclude and highlight directions for future work.

2 | NOTATION AND PRELIMINARIES

We review basic operations on matrices and tensors invoked throughout this paper, adopting the terminology and notation in Kolda and Bader (2009). Throughout the paper, we use lowercase letters to indicate scalars, bold lowercase letters to indicate vectors, bold capital characters to indicate matrices, and bold calligraphic capital characters to indicate tensors. We will also use the shorthand $[n]$ to denote an index set $\{1, \dots, n\}$.

Tensors can be considered generalizations of scalars, vectors, and matrices. Let \mathcal{X} represent an D -dimensional tensor in $\mathbb{R}^{p_1 \times \dots \times p_D}$. The tensor \mathcal{X} has order D , its number of dimensions or modes. For example, vectors are tensors of order one and have one mode. Matrices are tensors of order two and have two modes. We denote an element of \mathcal{X} by $x_{i_1 i_2, \dots, i_D}$, where $i_k \in [p_k]$ and $i \in [D]$. Fibres are the generalization of matrix rows and columns to higher order tensors. A fibre is defined by fixing the index of every dimension except one dimension. Mode i fibres are p_i -dimensional vectors extracted from \mathcal{X} by fixing all the indices $(i_1, \dots, i_{l-1}, i_{l+1}, \dots, i_D)$ except the i th one i_l . For example, columns of a matrix are Mode 1 fibres, and rows of a matrix are Mode 2 fibres.

It is often useful to reshape a tensor into a matrix. Reordering a tensor into a matrix is referred to as matricization. The mode i matricization of a tensor $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_D}$, denoted $\mathbf{X}_{(i)} \in \mathbb{R}^{p_i \times p_{-i}}$ with $p_{-i} = \prod_{k=1, k \neq i}^D p_k$, arranges the mode i fibres as the columns of the matrix $\mathbf{X}_{(i)}$. In a mode i matricization, the tensor element x_{i_1, \dots, i_D} is mapped to the matrix element of $\mathbf{X}_{(i)}$, with index (i, j) , where $j = 1 + \sum_{k=1, k \neq i}^D (i_k - 1)J_k$ with

$$J_k = \begin{cases} 1 & \text{if } k = 1 \text{ or if } k = 2 \text{ and } i = 1. \\ \prod_{k=1, k' \neq i}^{k-1} p_{k'} & \text{otherwise.} \end{cases}$$

Reordering a tensor into a vector is referred to as vectorization. We first describe vectorization of a matrix before describing vectorization of a general tensor. The vectorization of a matrix \mathbf{X} is denoted by $\text{vec}(\mathbf{X})$ and is the vector obtained by stacking the columns of \mathbf{X} on top of each other. The vectorization of the mode i matricization of a tensor \mathcal{X} in turn is denoted as $\mathbf{x}_{(i)} = \text{vec}(\mathbf{X}_{(i)})$. We then define the vectorization of a tensor \mathcal{X} , denoted by $\text{vec}(\mathcal{X})$, as the vectorization of its Mode 1 matricization, namely, $\text{vec}(\mathbf{X}_{(1)})$. When unambiguous from context, we will often denote the vectorization of a tensor \mathcal{X} by its corresponding bold lowercase \mathbf{x} .

The inner product of two tensors of compatible dimensions $\mathcal{X}, \tilde{\mathcal{X}} \in \mathbb{R}^{p_1 \times \dots \times p_D}$ is the sum of the product of their entries, namely,

$$\langle \mathcal{X}, \tilde{\mathcal{X}} \rangle = \sum_{i_1=1}^{p_1} \dots \sum_{i_D=1}^{p_D} x_{i_1, \dots, i_D} \tilde{x}_{i_1, \dots, i_D}.$$

The mode i product of a tensor $\tilde{\mathcal{X}} \in \mathbb{R}^{p_1 \times \dots \times p_D}$ with a matrix $\mathbf{U} \in \mathbb{R}^{J \times p_i}$ is denoted by $\mathcal{X} \times_i \mathbf{U}$ and is the tensor of size $p_1 \times \dots \times p_{i-1} \times J \times p_{i+1} \times \dots \times p_D$ with elements

$$(\mathcal{X} \times_i \mathbf{U})_{i_1, \dots, i_{i-1}, j, i_{i+1}, \dots, i_D} = \sum_{i_i=1}^{p_i} x_{i_1, i_2, \dots, i_D} U_{j i_i}.$$

Finally, we review three kinds of matrix products, as well as one definition of matrix division, that will be used throughout the paper.

- For two matrices $\mathbf{A} \in \mathbb{R}^{p_1 \times p_2}$ and $\mathbf{B} \in \mathbb{R}^{q_1 \times q_2}$, the Kronecker product is the $p_1 q_1$ -by- $p_2 q_2$ matrix,

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11} \mathbf{B} & \dots & a_{1 p_2} \mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{p_1 1} \mathbf{B} & \dots & a_{p_1 p_2} \mathbf{B} \end{pmatrix}.$$

- For two matrices $\mathbf{A} = (a_1 \dots a_{p_2}) \in \mathbb{R}^{p_1 \times p_2}$ and $\mathbf{B} = (b_1 \dots b_{p_2}) \in \mathbb{R}^{q_1 \times p_2}$ that have the same number of columns p_2 , the Khatri-Rao product is the $p_1 q_1$ -by- p_2 matrix,

$$\mathbf{A} \circ \mathbf{B} = (a_1 \otimes b_1 \ a_2 \otimes b_2 \ \dots \ a_{p_2} \otimes b_{p_2}).$$

which is a column-wise Kronecker product of \mathbf{A} and \mathbf{B} .

- For two matrices \mathbf{A} and \mathbf{B} of the same size, the Hadamard product is the element-wise product $\mathbf{A} * \mathbf{B} = \{a_{ij} b_{ij}\}$. Because the Hadamard product commutes, we use $*_i \mathbf{A}_i$ to denote $\mathbf{A}_1 * \dots * \mathbf{A}_m = \mathbf{A}_{\pi(1)} * \dots * \mathbf{A}_{\pi(m)}$ for any permutation π .
- Finally, for two matrices \mathbf{A} and \mathbf{B} of the same size, the Hadamard quotient is the element-wise quotient $\mathbf{A} \oslash \mathbf{B} = \{a_{ij}/b_{ij}\}$.

3 | TENSORCANONICAL CORRELATION ANALYSIS

3.1 | Population TCCA

Let $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_{D_x}}$ and $\mathcal{Y} \in \mathbb{R}^{q_1 \times \dots \times q_{D_y}}$ be two random tensors of order D_x and D_y , respectively. We denote the vectorizations of \mathcal{X} and \mathcal{Y} by \mathbf{x} and \mathbf{y} , respectively. Denote by Σ_x and Σ_y the covariances of \mathbf{x} and \mathbf{y} , respectively. Denote by $\Sigma_{x,y}$ the covariance between \mathbf{x}

and \mathbf{y} . Let $\mathcal{V} \in \mathbb{R}^{p_1 \times \dots \times p_{D_x}}$, and $\mathcal{W} \in \mathbb{R}^{q_1 \times \dots \times q_{D_y}}$ be constant tensors, and let $\rho(\mathcal{V}, \mathcal{W})$ denote the correlation between the two linear combinations $\langle \mathcal{X}, \mathcal{V} \rangle$ and $\langle \mathcal{Y}, \mathcal{W} \rangle$ namely,

$$\rho(\mathcal{V}, \mathcal{W}) = \frac{\text{Cov}(\langle \mathcal{X}, \mathcal{V} \rangle, \langle \mathcal{Y}, \mathcal{W} \rangle)}{\sqrt{\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle)}\sqrt{\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle)}} = \frac{\mathbf{v}^\top \Sigma_{\mathbf{x}, \mathbf{y}} \mathbf{w}}{\sqrt{\mathbf{v}^\top \Sigma_{\mathbf{x}} \mathbf{v}} \sqrt{\mathbf{w}^\top \Sigma_{\mathbf{y}} \mathbf{w}}}. \quad (1)$$

The pair $(\mathcal{V}, \mathcal{W})$ that maximizes ρ are the canonical tensors, and the optimal is the canonical coefficient. Maximizing the objective in Equation (1) presents two challenges: (a) high dimensionality of optimization variables \mathcal{V} and \mathcal{W} and (b) the estimation of the huge covariance matrices $\Sigma_{\mathbf{x}}$ and $\Sigma_{\mathbf{y}}$ and the cross-covariance matrix $\Sigma_{\mathbf{x}, \mathbf{y}}$. We will address challenge (b) by imposing a separable covariance structure in Section 4.

To address challenge (a), we impose the parsimonious CANDECOMP/PARAFAC (CP), or Kruskal, representation on the canonical tensors. The CP representation generalizes the idea of representing a matrix as the sum of Rank 1 matrices to representing a tensor as the sum of Rank 1 tensors. An order- D tensor $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_D}$ is Rank 1 if it can be expressed as the outer product of D vectors $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \dots, \mathbf{a}^{(D)}$, namely, $\mathcal{X} = \mathbf{a}^{(1)} \circ \mathbf{a}^{(2)} \circ \dots \circ \mathbf{a}^{(D)}$, where the binary operator \circ denotes the vector outer product. Thus, the $(\iota_1, \iota_2, \dots, \iota_D)$ th element of \mathcal{X} is $x_{\iota_1 \iota_2 \dots \iota_D} = a_{\iota_1}^{(1)} a_{\iota_2}^{(2)} \dots a_{\iota_D}^{(D)}$. A rank- R tensor can be written as the sum of R Rank 1 tensors, namely,

$$\mathcal{X} = [A_1, \dots, A_D] = \sum_{r=1}^R a_r^{(1)} \circ \dots \circ a_r^{(D)},$$

where $A_i = (a_1^{(i)} \ a_2^{(i)} \ \dots \ a_R^{(i)}) \in \mathbb{R}^{p_i \times R}$ denotes the mode i factor matrix. We use the Kruskal notation [...] to concisely summarize the sum.

Thus, instead of searching over the space of all order- D_x and order- D_y tensor pairs $(\mathcal{V}, \mathcal{W})$, we limit our search to tensors of rank- R_x and rank- R_y ,

$$\mathcal{V} = [V_1, \dots, V_{D_x}], V_i \in \mathbb{R}^{p_i \times R_x}, i \in [D_x], \mathcal{W} = [W_1, \dots, W_{D_y}], W_j \in \mathbb{R}^{q_j \times R_y}, j \in [D_y]. \quad (2)$$

As we will see later, this parameterization makes progress towards alleviating the burden of estimating a huge covariance matrix. Note that

$$\langle \mathcal{X}, \mathcal{V} \rangle = \langle \mathcal{X}, \sum_{r=1}^{R_x} \mathbf{w}_r^{(1)} \circ \dots \circ \mathbf{w}_r^{(D_x)} \rangle = \sum_{r=1}^{R_x} x \times_1 \mathbf{w}_r^{(1)} \times_2 \dots \times_{D_x} \mathbf{w}_r^{(D_x)}$$

$$\langle \mathcal{Y}, \mathcal{W} \rangle = \langle \mathcal{Y}, \sum_{r=1}^{R_Y} \mathbf{w}_r^{(1)} \circ \dots \circ \mathbf{w}_r^{(D_Y)} \rangle = \sum_{r=1}^{R_Y} \mathcal{Y} \times_1 \mathbf{w}_r^{(1)} \times_2 \dots \times_{D_Y} \mathbf{w}_r^{(D_Y)}.$$

Thus, we seek to maximize the correlation between a rank- R_X multilinear form in \mathcal{X} and a rank- R_Y multilinear form in \mathcal{Y} . Multiway information is preserved, and the dimensionality is reduced from an exponential number of parameters $\prod_{i=1}^{D_X} p_i + \prod_{j=1}^{D_Y} q_j$ to a linear number of parameters $R_X \sum_{i=1}^{D_X} p_i + R_Y \sum_{j=1}^{D_Y} q_j$. Note that the ranks (R_X, R_Y) here are not the number of canonical tensor pairs being sought. In this paper, we focus on obtaining only the top canonical tensor pair $(\mathcal{V}, \mathcal{W})$, which have ranks R_X and R_Y .

The following representations of $\text{Var}(\langle \mathcal{X}, \mathcal{Y} \rangle)$, $\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle)$, and $\text{Cov}(\langle \mathcal{X}, \mathcal{Y} \rangle, \langle \mathcal{Y}, \mathcal{W} \rangle)$ in terms of a CP decomposition are keys to our estimation algorithms. Its proof is in the Supporting Information.

Proposition 1—Let $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_{D_X}}$ and $\mathcal{Y} \in \mathbb{R}^{q_1 \times \dots \times q_{D_Y}}$ be two random tensors and $\mathcal{V} = [\mathbf{V}_1, \dots, \mathbf{V}_{D_X}]$ and $\mathcal{W} = [\mathbf{W}_1, \dots, \mathbf{W}_{D_Y}]$ be two constant tensors of the same size as \mathcal{X} and \mathcal{Y} respectively. Define

$$\mathbf{V}_{(-i)} = [\mathbf{V}_{D_X} \circ \dots \circ \mathbf{V}_{i+1} \circ \mathbf{V}_{i-1} \circ \dots \circ \mathbf{V}_1] \otimes \mathbf{I}_{p_i}, \quad (3)$$

$$\mathbf{W}_{(-j)} = [\mathbf{W}_{D_Y} \circ \dots \circ \mathbf{W}_{j+1} \circ \mathbf{W}_{j-1} \circ \dots \circ \mathbf{W}_1] \otimes \mathbf{I}_{q_j}, \quad (4)$$

and let $\Sigma_{\mathbf{x}_{(i)}}$ denote the covariance of $\mathbf{x}_{(i)} = \text{vec}(\mathbf{X}_{(i)})$. Define $\Sigma_{\mathbf{y}_{(j)}}$ and $\Sigma_{\mathbf{x}_{(i)}, \mathbf{y}_{(j)}}$ analogously.

Then

$$\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle) = \mathbf{v}_i^\top \mathbf{V}_{(-i)}^\top \Sigma_{\mathbf{x}_{(i)}} \mathbf{V}_{(-i)} \mathbf{v}_i$$

$$\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle) = \mathbf{w}_j^\top \mathbf{W}_{(-j)}^\top \Sigma_{\mathbf{y}_{(j)}} \mathbf{W}_{(-j)} \mathbf{w}_j$$

$$\text{Cov}(\langle \mathcal{X}, \mathcal{V} \rangle, \langle \mathcal{Y}, \mathcal{W} \rangle) = \mathbf{v}_i^\top \mathbf{V}_{(-i)}^\top \Sigma_{\mathbf{x}_{(i)}, \mathbf{y}_{(j)}} \mathbf{W}_{(-j)} \mathbf{w}_j$$

for any $i \in [D_X]$ and $j \in [D_Y]$, where $\mathbf{v}_i = \text{vec}(\mathbf{V}_i)$ and $\mathbf{w}_j = \text{vec}(\mathbf{W}_j)$.

3.2 | Sample TCCA

Suppose we observe N pairs of i.i.d. tensor data $(\mathcal{X}_n, \mathcal{Y}_n)$, and we estimate \mathcal{V} and \mathcal{W} by solving the optimization problem

$$\text{maximize } \hat{\rho}(\mathcal{V}, \mathcal{W}) \equiv \frac{\mathbf{v}^\top \hat{\Sigma}_{\mathbf{x}, \mathbf{y}} \mathbf{w}}{\sqrt{\mathbf{v}^\top \hat{\Sigma}_{\mathbf{x}} \mathbf{v} \mathbf{w}^\top \hat{\Sigma}_{\mathbf{y}} \mathbf{w}}}. \quad (5)$$

where $\hat{\Sigma}_{\mathbf{x}}$, $\hat{\Sigma}_{\mathbf{y}}$ and $\hat{\Sigma}_{\mathbf{x}, \mathbf{y}}$ are sample estimates of the corresponding covariances. Recall that CCA models can be estimated numerically by computing the solution to the following generalized eigenvalue problem:

$$\begin{pmatrix} 0 & \hat{\Sigma}_{\mathbf{x}, \mathbf{y}} \\ \hat{\Sigma}_{\mathbf{y}, \mathbf{x}} & 0 \end{pmatrix} \begin{pmatrix} v \\ w \end{pmatrix} = \rho \begin{pmatrix} \hat{\Sigma}_{\mathbf{x}} & 0 \\ 0 & \hat{\Sigma}_{\mathbf{y}} \end{pmatrix} \begin{pmatrix} v \\ w \end{pmatrix}.$$

This problem is guaranteed to have a solution if and only if the covariance matrices $\hat{\Sigma}_{\mathbf{x}}$ and $\hat{\Sigma}_{\mathbf{y}}$ are nonsingular. In practice, the sample size N is smaller than the size of $\hat{\Sigma}_{\mathbf{x}}$ and $\hat{\Sigma}_{\mathbf{y}}$, $\left(\prod_{i=1}^{D_x} p_i\right) \left(\prod_{i=1}^{D_x} p_i\right)$, $\left(\prod_{j=1}^{D_y} q_j\right) \times \left(\prod_{j=1}^{D_y} q_j\right)$, respectively. Therefore, the sample covariance matrices $\hat{\Sigma}_{\mathbf{x}}$ and $\hat{\Sigma}_{\mathbf{y}}$ are singular, and a solution cannot be obtained. As a remedy, several regularized estimation methods for obtaining nonsingular sample covariance matrices (Vinod, 1976; Ledoit & Wolf, 2004; González, Déjean, Martin, & Baccini, 2008; Ledoit & Wolf, 2012; Cai & Yuan, 2012; Bickel & Levina, 2008; 2008; Kubokawa et al., 2013; Srivastava & Reid, 2012) can be used.

If we take $\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle)$, $\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle)$ and $\text{Cov}(\langle \mathcal{X}, \mathcal{V} \rangle, \langle \mathcal{Y}, \mathcal{W} \rangle)$ to be $\mathbf{v}^\top \hat{\Sigma}_{\mathbf{x}} \mathbf{v}$, $\mathbf{w}^\top \hat{\Sigma}_{\mathbf{y}} \mathbf{w}$, and $\mathbf{v}^\top \hat{\Sigma}_{\mathbf{x}, \mathbf{y}} \mathbf{w}$ respectively, then Proposition 1 suggests a block coordinate ascent algorithm where we update the factor matrices *in pairs* $(\mathbf{V}_i, \mathbf{W}_j)$, for different combinations of $i \in [D_x]$ and $j \in [D_y]$. To update the pair $(\mathbf{V}_i, \mathbf{W}_j)$, we solve the following problem:

$$\begin{aligned} & \text{maximize } \mathbf{v}_i^\top \mathbf{V}_{(-i)}^\top \hat{\Sigma}_{\mathbf{x}(i), \mathbf{y}(j)} \mathbf{W}_{(-j)} \mathbf{w}_j \text{ subject to } \mathbf{v}_i^\top \mathbf{V}_{(-i)}^\top \hat{\Sigma}_{\mathbf{x}(i)} \mathbf{V}_{(-i)} \mathbf{v}_i = 1, \\ & \mathbf{w}_j^\top \mathbf{W}_{(-j)}^\top \hat{\Sigma}_{\mathbf{y}(j)} \mathbf{W}_{(-j)} \mathbf{w}_j = 1, \end{aligned} \quad (6)$$

which is a substantially smaller optimization problem over $p_i R_x + q_j R_y$ variables compared with any alternative “all-at-once” strategies to iteratively optimize over all $\left(\prod_{i=1}^{D_x} p_i\right) \left(\prod_{i=1}^{D_x} p_i\right) + \left(\prod_{j=1}^{D_y} q_j\right) \times \left(\prod_{j=1}^{D_y} q_j\right)$ parameters simultaneously. Problem (6) includes $\hat{\Sigma}_{\mathbf{x}(i)}$, a permuted version of $\hat{\Sigma}_{\mathbf{x}}$ with size $\left(\prod_{i=1}^{D_x} p_i\right) \left(\prod_{i=1}^{D_x} p_i\right)$. However, we work on the “compressed” covariance matrix

$$\mathbf{V}_{(-i)}^\top \hat{\Sigma}_{\mathbf{x}(i)} \mathbf{V}_{(-i)} \in \mathbb{R}^{p_i R_x \times p_i R_x},$$

which is likely to be full rank when $N > p_i R_x$ instead of a singular matrix $\hat{\Sigma}_{\mathbf{x}(i)}$. This approach enables us to solve the generalized eigenvalue problem with matrices that are $(p_i R_x) + (q_j R_y)$ -by- $(p_i R_x) + (q_j R_y)$. Algorithm 1 summarizes the estimation procedure that

comes with the following convergence guarantees. The proof is in the Supporting Information.

Proposition 2—If the matrices $\widehat{\Sigma}_{\mathbf{x}}$ and $\widehat{\Sigma}_{\mathbf{y}}$ are nonsingular, then the limit Points of the iterate sequence generated by Algorithm 1 are canonical tensors of the sample TCCA Problem.

Algorithm 1

TCCA, with (i) assumptions on CP structure on canonical tensors (\mathcal{V} , \mathcal{W}) and (ii) no additional assumptions on the structures of the covariances $\text{Var}(\text{vec}(\mathcal{X}))$ and $\text{Var}(\text{vec}(\mathcal{Y}))$

Initialize $v_i^{(0)}$ and $w_j^{(0)}$, for $i \in [D_x]$ and $j \in [D_y]$

$t \leftarrow 0$

repeat

Select $(i, j) \in [D_x] \times [D_y]$

$\mathbf{V}_{(-i)}^{(t)} \leftarrow \left[\mathbf{v}_{D_x}^{(t)} \odot \cdots \odot \mathbf{v}_{i+1}^{(t)} \odot \mathbf{v}_{i-1}^{(t)} \odot \cdots \odot \mathbf{v}_1^{(t)} \right] \otimes \mathbf{I}_{p_i}$

$\mathbf{W}_{(-j)}^{(t)} \leftarrow \left[\mathbf{w}_{D_y}^{(t)} \odot \cdots \odot \mathbf{w}_{j+1}^{(t)} \odot \mathbf{w}_{j-1}^{(t)} \odot \cdots \odot \mathbf{w}_1^{(t)} \right] \otimes \mathbf{I}_{q_j}$

$\mathbf{C}_{\mathbf{x}}^{(t)} \leftarrow \mathbf{V}_{(-i)}^{(t)T} \widehat{\Sigma}_{\mathbf{x}(i)} \mathbf{V}_{(-i)}^{(t)}$

$\mathbf{C}_{\mathbf{y}}^{(t)} \leftarrow \mathbf{W}_{(-j)}^{(t)\top} \widehat{\Sigma}_{\mathbf{y}(j)} \mathbf{W}_{(-j)}^{(t)}$

$\mathbf{C}_{\mathbf{x}, \mathbf{y}}^{(t)} \leftarrow \mathbf{V}_{(-i)}^{(t)T} \widehat{\Sigma}_{\mathbf{x}(i)\mathbf{y}(j)} \mathbf{W}_{(-j)}^{(t)}$

Generalized eigen-decomposition: $\begin{pmatrix} 0 & \mathbf{C}_{\mathbf{x}, \mathbf{y}}^{(t)} \\ \mathbf{C}_{\mathbf{x}, \mathbf{y}}^{(t)T} & 0 \end{pmatrix} \begin{pmatrix} w_i^{(t+1)} \\ w_j^{(t+1)} \end{pmatrix} = \rho^{(t+1)} \begin{pmatrix} \mathbf{C}_{\mathbf{x}}^{(t)} & 0 \\ 0 & \mathbf{C}_{\mathbf{y}}^{(t)} \end{pmatrix} \begin{pmatrix} w_i^{(t+1)} \\ w_j^{(t+1)} \end{pmatrix}$

$t \leftarrow t+1$

until $\rho^{(t)}$ converges

4 | TCCA WITH SEPARABLE COVARIANCE STRUCTURE

4.1 | Population TCCA with separable covariances

Hoff (2011) proposed the array normal distribution with separable covariance structure. Separable marginal covariances are defined as

$$\Sigma_{\mathbf{x}} = \text{Var}(\mathbf{x}) = \Sigma_{\mathbf{x}, D_x} \otimes \cdots \otimes \Sigma_{\mathbf{x}, 1} \text{ and } \Sigma_{\mathbf{y}} = \text{Var}(\mathbf{y}) = \Sigma_{\mathbf{y}, D_y} \otimes \cdots \otimes \Sigma_{\mathbf{y}, 1}, \quad (7)$$

where $\Sigma_{\mathbf{x}, j} \in \mathbb{R}^{p_i \times p_i}$ for $i \in [D_x]$ and $\Sigma_{\mathbf{y}, j} \in \mathbb{R}^{q_j \times q_j}$ for $j \in [D_y]$. Then the overall covariance of population model is

$$\text{Var} \begin{pmatrix} \mathbf{X} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \Sigma_{\mathbf{x}, D_x} \otimes \cdots \otimes \Sigma_{\mathbf{x}, 1} & \Sigma_{\mathbf{x}, \mathbf{y}} \\ \Sigma_{\mathbf{y}, \mathbf{x}} & \Sigma_{\mathbf{y}, D_y} \otimes \cdots \otimes \Sigma_{\mathbf{y}, 1} \end{pmatrix}. \quad (8)$$

Intuitively, $\Sigma_{x,i}$ summarizes the covariance along the mode i fibres of tensor \mathcal{X} and $\Sigma_{y,j}$ summarizes the covariance along the mode j fibres of tensor \mathcal{Y} . The following result shows a representation of $\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle)$ and $\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle)$ in presence of separable covariance structure (7) that we will leverage in our estimation algorithm.

Proposition 3—Let $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_{D_x}}$ and $\mathcal{Y} \in \mathbb{R}^{q_1 \times \dots \times q_{D_y}}$ be two random tensors admitting the separable covariance structure (7) and $\mathcal{V} = [v_1, \dots, v_{D_x}]$ and $\mathcal{W} = [W_1, \dots, W_{D_y}]$ be two constant tensors. Define

$$H_{x,-i} = *_{i' \neq i} (V_{i'}^\top \Sigma_{x,i'} V_{i'}) \text{ and } H_{y,-j} = *_{j' \neq j} (W_{j'}^\top \Sigma_{y,j'} W_{j'}).$$

Then

$$\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle) = v_i^\top (H_{x,-i} \otimes \Sigma_{x,i}) v_i \text{ and } \text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle) = w_j^\top (H_{y,-j} \otimes \Sigma_{y,j}) w_j.$$

for any $i \in [D_x]$ and $j \in [D_y]$, where $v_i = \text{vec}(V_i)$ and $w_j = \text{vec}(W_j)$.

With the separable covariance structure and the CP structure on canonical tensors (\mathcal{V}, \mathcal{W}), the objective function of TCCA population model (1) greatly simplifies. Note that the separable covariance structure may not hold for real data, in which case the covariance estimates are biased. Despite this drawback, this parsimonious structure is worth considering due to the stability that it can impart by reducing estimation variance.

4.2 | Sample TCCA with separable covariances

Given data $(\mathcal{X}_n, \mathcal{Y}_n)$, $n \in [N]$, the goal is to maximize the sample canonical correlation (5) under assumptions that (a) \mathcal{V} and \mathcal{W} have the CP decomposition structure and (b) \mathcal{X} and \mathcal{Y} admit the separable covariance structure. We follow the same strategy as sample TCCA in Section 3.2, updating parameters in pairs of factor matrices (V_i, W_j) . To update the pair (V_i, W_j) , we solve the subproblem

$$\text{maximize } v_i^\top V_{(-i)}^\top \widehat{\Sigma}_{x(i), y(j)} W_{(-j)} w_j, \text{ subject to } v_i^\top (H_{x,-i} \otimes \widehat{\Sigma}_{x,i}) v_i = 1, w_j^\top (H_{y,-j} \otimes \widehat{\Sigma}_{y,j}) w_j = 1.$$

where $H_{x,-i}$ and $H_{y,-j}$, defined in Proposition 3, are evaluated at current iterates $V_{i'}$ where $i' \neq i$, and $W_{j'}$, where $j' \neq j$.

By assuming the separable covariance structure, Proposition 3 enables us to greatly simplify the variance calculations for $\text{Var}(\langle \mathcal{X}, \mathcal{V} \rangle)$ and $\text{Var}(\langle \mathcal{Y}, \mathcal{W} \rangle)$. Note that the calculation of $V_{(-i)}^\top \widehat{\Sigma}_{x_0} V_{(-i)}$ and $W_{(-j)}^\top \widehat{\Sigma}_{y(j)} W_{(-j)}$ in the subproblem of the sample TCCA costs $(R_x \prod_{i=1}^{D_x} p_i)^2 + (R_y \prod_{j=1}^{D_y} q_j)^2$ flops. In contrast, the calculation of matrices $H_{x,-i} \otimes \widehat{\Sigma}_{x,i}$ and $H_{y,-j} \otimes \widehat{\Sigma}_{y,j}$ in the subproblem of the sample TCCA with the separable covariance structure only costs $(R_x p_i)^2 + (R_y q_j)^2$ flops. Algorithm 2 summarizes the estimation

procedure under the separable covariance assumption (7). Like Algorithm 1, Algorithm 2 also comes with convergence guarantees. The proof is in the Supporting Information.

Proposition 4—If the matrices $\widehat{\Sigma}_{\mathbf{x}}$ and $\widehat{\Sigma}_{\mathbf{y}}$ are nonsingular, then the limit points of the iterate sequence generated by Algorithm 2 are canonical tensors of the sample TCCA problem with separable covariances.

Algorithm 2

TCCA for two tensors of modes D_x and D_y , respectively, assuming (i) CP structure on canonical correlation tensors (\mathcal{V} , \mathcal{W}) and (ii) separable covariances $\text{Var}(\text{vec}(\mathcal{X}))$ and $\text{Var}(\text{vec}(\mathcal{Y}))$

Initialize $\mathbf{V}_i^{(0)}, i \in [D_x]$

$\mathbf{W}_j^{(0)}, j \in [D_y]$

$\mathbf{H}_{\mathbf{x}}^{(0)} \leftarrow *_i (\mathbf{V}_i^{(0)T} \widehat{\Sigma}_{\mathbf{x},i} \mathbf{V}_i^{(0)})$

$\mathbf{H}_{\mathbf{y}}^{(0)} \leftarrow *_j (\mathbf{W}_j^{(0)T} \widehat{\Sigma}_{\mathbf{y},j} \mathbf{W}_j^{(0)})$

$t \leftarrow 0$

repeat

Select $(i, j) \in [D_x] \times [D_y]$

$\mathbf{H}_{\mathbf{x},-i}^{(t)} \leftarrow \mathbf{H}_{\mathbf{x}}^{(t)} \tau(\mathbf{V}_i^{(t)T} \widehat{\Sigma}_{\mathbf{x},i} \mathbf{V}_i^{(t)})$

$\mathbf{C}_{\mathbf{x}}^{(t)} \leftarrow \mathbf{H}_{\mathbf{x},-i}^{(t)} \otimes \widehat{\Sigma}_{\mathbf{x},i}$

$\mathbf{H}_{\mathbf{y},-j}^{(t)} \leftarrow \mathbf{H}_{\mathbf{y}}^{(t)} \tau(\mathbf{W}_j^{(t)T} \widehat{\Sigma}_{\mathbf{y},j} \mathbf{W}_j^{(t)})$

$\mathbf{C}_{\mathbf{y}}^{(t)} \leftarrow \mathbf{H}_{\mathbf{y},-j}^{(t)} \otimes \widehat{\Sigma}_{\mathbf{y},j}$

$\mathbf{V}_{(-i)}^{(t)} \leftarrow \mathbf{V}_{D_x}^{(t)} \odot \dots \odot \mathbf{V}_{i+1}^{(t)} \odot \mathbf{V}_{i-1}^{(t)} \odot \dots \odot \mathbf{V}_1^{(t)}$

$\mathbf{W}_{(-j)}^{(t)} \leftarrow \mathbf{W}_{D_y}^{(t)} \odot \dots \odot \mathbf{W}_{j+1}^{(t)} \odot \mathbf{W}_{j-1}^{(t)} \odot \dots \odot \mathbf{W}_1^{(t)}$

$\mathbf{C}_{\mathbf{xy}}^{(t)} \leftarrow \left(\mathbf{V}_{-i}^{(t)T} \otimes \mathbf{I}_{p_i} \right) \widehat{\Sigma}_{\mathbf{x}(i),\mathbf{y}(j)} \left(\mathbf{W}_{-j}^{(t)T} \otimes \mathbf{I}_{q_j} \right)$

Solve following generalized eigenvalue decomposition, $\begin{bmatrix} \mathbf{0} & \mathbf{C}_{\mathbf{xy}}^{(t)} \\ \mathbf{C}_{\mathbf{xy}}^{(t)} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{v}_i^{(t+1)} \\ \mathbf{w}_j^{(t+1)} \end{bmatrix} = \rho^{(t+1)} \begin{bmatrix} \mathbf{C}_{\mathbf{x}}^{(t)} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{\mathbf{y}}^{(t)} \end{bmatrix} \begin{bmatrix} \mathbf{v}_i^{(t+1)} \\ \mathbf{w}_j^{(t+1)} \end{bmatrix}$

$\mathbf{H}_{\mathbf{x}}^{(t+1)} \leftarrow \mathbf{H}_{\mathbf{x},-i}^{(t)} *_i \left(\mathbf{v}_i^{(t+1)T} \widehat{\Sigma}_{\mathbf{x},i} \mathbf{v}_i^{(t+1)} \right)$

$\mathbf{H}_{\mathbf{y}}^{(t+1)} \leftarrow \mathbf{H}_{\mathbf{y},-j}^{(t)} *_j \left(\mathbf{w}_j^{(t+1)T} \widehat{\Sigma}_{\mathbf{y},j} \mathbf{w}_j^{(t+1)} \right)$

$t \leftarrow t+1$

until $\rho^{(t)}$ converges

4.2.1 | Estimation of separable covariance matrices—Algorithm 2 relies on the sample estimate of separable covariances $\widehat{\Sigma}_{x,i}$ and $\widehat{\Sigma}_{y,j}$, and the unstructured cross-covariance $\widehat{\Sigma}_{x,y}$. The following lemma will be useful in computing a consistent estimator for covariance matrices with the separable structure.

Lemma 1: If a random tensor $\mathcal{X} \in \mathbb{R}^{p_1 \times \dots \times p_D}$ has mean zero and separable covariance $\text{Var}(\mathbf{x}) = \Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1}$, then

$$\mathbb{E}(\mathbf{X}_{(i)}\mathbf{X}_{(i)}^\top) = \left(\prod_{i' \neq i} \text{tr}(\Sigma_{x,i'}) \right) \Sigma_{x,i} \text{ and } \mathbb{E}(\|\mathbf{x}\|_2^2) = \prod_{i=1}^{D_x} \text{tr}(\Sigma_{x,i}).$$

Proof. For the first identity, see Hoff, 2011 (2011, Proposition 2.1). For the second identity, $\mathbb{E}(\|\mathbf{x}\|_2^2) = \text{tr}(\text{Var}(\mathbf{x})) = \text{tr}(\Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1}) = \prod_i \text{tr}(\Sigma_{x,i})$. \square

Given N i.i.d. observations $(\mathbf{x}_n, \mathbf{y}_n)$, the estimators $\widehat{r}_x = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \bar{\mathbf{x}}\|_2^2$ and $\widehat{r}_y = \frac{1}{N} \sum_{n=1}^N \|\mathbf{y}_n - \bar{\mathbf{y}}\|_2^2$ consistently estimate $\prod_{i=1}^{D_x} \text{tr}(\Sigma_{x,i})$ and $\prod_{j=1}^{D_y} \text{tr}(\Sigma_{y,j})$, respectively, where $\bar{\mathbf{x}}$ and $\bar{\mathbf{y}}$ are the sample means of the vectorized tensors \mathbf{x}_n and \mathbf{y}_n . We propose the following covariance estimators:

$$\begin{aligned} \widehat{\Sigma}_{x,i} &= \frac{1}{N r_x^{(D_x-1)/D_x}} \sum_{n=1}^N (x_{n(i)} - \bar{x}_{(i)})(x_{n(i)} - \bar{x}_{(i)})^\top, \widehat{\Sigma}_{y,j} \\ &= \frac{1}{N \widehat{r}_y^{(D_y-1)/D_y}} \sum_{n=1}^N (y_{n(j)} - \bar{y}_{(j)})(y_{n(j)} - \bar{y}_{(j)})^\top, \end{aligned}$$

$$\widehat{\Sigma}_{x,y} = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{\mathbf{x}})(y_n - \bar{\mathbf{y}})^\top,$$

where $i \in [D_x]$ and $j \in [D_y]$. The vectors $\mathbf{x}_{n(i)}$ and $\bar{x}_{(i)}$ denote the mode i vectorization of the n th observation and the sample mean of the mode i vectorized tensors \mathbf{x}_n , respectively. The vectors $\mathbf{y}_{n(j)}$ and $\bar{y}_{(j)}$ denote analogous vectorizations.

Unfortunately, the separable covariance structure (8) is not identifiable in the individual $\Sigma_{x,i}$ due to scaling indeterminacy. Therefore, $\Sigma_{x,i}$ cannot be consistently estimated. Note, however, that we do not need to consistently estimate the individual $\Sigma_{x,i}$ in order to consistently estimate their Kronecker product. To see this, note that by Slutsky's theorem,

$$\widehat{\Sigma}_{x,D_x} \otimes \dots \otimes \widehat{\Sigma}_{x,1} \rightarrow \frac{\left(\prod_i \text{tr}(\Sigma_{x,i}) \right)^{D_x-1}}{\left(\prod_i \text{tr}(\Sigma_{x,i}) \right)^{D_x-1}} \Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1} = \Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1}$$

consistently estimates $\text{Var}(\mathbf{x})$.

Note that Hoff (2011) proposes an iterative algorithm for finding the maximum likelihood estimation (MLE) on the basis of the array normal assumption. The maximum likelihood estimation may improve upon the above estimates when data actually come from an array normal distribution.

5 | SPARSE TCCA

We may also recover sparse canonical tensors for both TCCA and scTCCA to enhance the interpretability of the estimated canonical tensors. Following the iterative thresholding strategy introduced by Ma (2013) for sparse principal component analysis, by Wang, Gu, Ning, and Liu (2015) for sparse expectation-maximization algorithms, and by Tan et al. (2018) for generalized eigenvalue problems, we incorporate a hard-thresholding step in Algorithms 1 and 2, as follows:

$$\mathbf{w}_i^{(t+1)} \leftarrow \Theta_\lambda(\mathbf{w}_i^{(t+1)}) \text{ and } \mathbf{w}_j^{(t+1)} \leftarrow \Theta_\lambda(\mathbf{w}_j^{(t+1)}),$$

where $\Theta_\lambda(\mathbf{v})$ performs element-wise hard-thresholding on \mathbf{v} , namely, the i th element of $\Theta_\lambda(\mathbf{v})$ is v_i if $|v_i| > \lambda$ and 0 otherwise. In the simulation studies of Section 7, we employ fivefold cross-validation with a grid point search to choose the tuning parameter λ , following Tan, Wang, Liu, and Zhang (2018). Instead of a grid point search, a random search would be another choice (Bergstra & Bengio, 2012).

6 | PROBABILISTIC MODEL FOR TCCA

Bach and Jordan (2006) give a probabilistic interpretation for the classic CCA, which enables us to simulate data with desired canonical vectors and canonical correlations. For TCCA without any assumption on the covariance structure, it is the same as the regular CCA by treating vectorized canonical tensors as canonical vectors. In this section, we first discuss how to generate data from given d canonical correlations $\rho_d = (\rho_1, \dots, \rho_d)$ and their corresponding canonical vectors in columns of matrices $(\mathbf{V}_d, \mathbf{W}_d)$, and then we extend it to TCCA with the separable covariance assumption. Let Σ_x and Σ_y be two covariance matrices such that $\mathbf{V}_d^\top \Sigma_x \mathbf{V}_d = \mathbf{W}_d^\top \Sigma_y \mathbf{W}_d = \mathbf{I}_d$. Define two linear transformations $A_x = \Sigma_x \mathbf{V}_d \mathbf{M}_x$ and $A_y = \Sigma_y \mathbf{W}_d \mathbf{M}_y$, where $\mathbf{M}_x, \mathbf{M}_y \in \mathbb{R}^{d \times d}$ are arbitrary matrices such that $\mathbf{M}_x \mathbf{M}_y^\top = \text{diag}(\rho_d)$ and their spectral norms are less than 1. We consider the latent factor model

$$\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_d)$$

$$\mathbf{x} | \mathbf{Z} = \mathbf{z} \sim N(\mathbf{A}_x \mathbf{z} + \boldsymbol{\mu}_x, \Sigma_x - \mathbf{A}_x \mathbf{A}_x^\top)$$

$$\mathbf{y} \mid \mathbf{Z} = \mathbf{z} \sim N(\mathbf{A}_y \mathbf{z} + \boldsymbol{\mu}_y, \Sigma_y - \mathbf{A}_y \mathbf{A}_y^\top).$$

The joint distribution of (\mathbf{x}, \mathbf{y}) is

$$\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \sim N \left(\begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}, \begin{bmatrix} \Sigma_x \Sigma_{xy} \\ \Sigma_{xy}^\top \Sigma_y \end{bmatrix} \right),$$

where $\Sigma_{xy} = \Sigma_x \mathbf{V}_d \text{diag}(\rho_d) \mathbf{W}_d^\top \Sigma_y$.

Now, we discuss how to construct Σ_x and Σ_y from $(\mathbf{V}_d, \mathbf{W}_d)$, which is not described in Bach and Jordan (2006). Let $\mathbf{V}_d = \mathbf{Q}_x \mathbf{R}_x$ be the thin QR decomposition of \mathbf{V}_d . Then

$$\Sigma_x = \mathbf{Q}_x \mathbf{R}_x^{-T} \mathbf{R}_x^{-1} \mathbf{Q}_x^\top + \mathbf{T}_x (\mathbf{I}_p - \mathbf{Q}_x \mathbf{Q}_x^\top) \mathbf{T}_x^\top$$

satisfies $\mathbf{V}_d^\top \Sigma_x \mathbf{V}_d = \mathbf{I}_d$ for arbitrary $\mathbf{T}_x \in \mathbb{R}^{p \times p}$. Similarly, let $\mathbf{W}_d = \mathbf{Q}_y \mathbf{R}_y$ be the thin QR decomposition of \mathbf{W}_d . Then

$$\Sigma_y = \mathbf{Q}_y \mathbf{R}_y^{-T} \mathbf{R}_y^{-1} \mathbf{Q}_y^\top + \mathbf{T}_y (\mathbf{I}_q - \mathbf{Q}_y \mathbf{Q}_y^\top) \mathbf{T}_y^\top$$

satisfies $\mathbf{W}_d^\top \Sigma_y \mathbf{W}_d = \mathbf{I}_d$ for arbitrary $\mathbf{T}_y \in \mathbb{R}^{q \times q}$. In this notation, the joint covariance is

$$\begin{aligned} \text{Var} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} &= \begin{pmatrix} \mathbf{Q}_x \mathbf{R}_x^{-T} \mathbf{R}_x^{-1} \mathbf{Q}_x^\top & \mathbf{Q}_x \mathbf{R}_x^{-T} \text{diag}(\rho_d) \mathbf{R}_y^\top \mathbf{Q}_y^\top \\ \mathbf{Q}_y \mathbf{R}_y^{-T} \text{diag}(\rho_d) \mathbf{R}_y^\top \mathbf{Q}_y^\top & \mathbf{Q}_x \mathbf{R}_x^{-T} \mathbf{R}_y^\top \mathbf{Q}_y^\top \end{pmatrix} \\ &+ \begin{pmatrix} \mathbf{T}_x (\mathbf{I}_p - \mathbf{Q}_x \mathbf{Q}_x^\top) \mathbf{T}_x^\top & \mathbf{0} \\ \mathbf{0} & \mathbf{T}_y (\mathbf{I}_q - \mathbf{Q}_y \mathbf{Q}_y^\top) \mathbf{T}_y^\top \end{pmatrix}. \end{aligned} \tag{9}$$

\mathbf{T}_x and \mathbf{T}_y are free parameters that adjust the noise level in \mathbf{x} and \mathbf{y} , respectively. The normal generative model is

$$\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_d)$$

$$\mathbf{x} \mid \mathbf{z} = \mathbf{z} \sim N(\mathbf{Q}_x \mathbf{R}_x^{-T} \mathbf{M}_x \mathbf{z} + \boldsymbol{\mu}_x, \Sigma_x)$$

$$\mathbf{y} \mid \mathbf{z} = \mathbf{z} \sim N(\mathbf{Q}_y \mathbf{R}_y^{-T} \mathbf{M}_y \mathbf{z} + \boldsymbol{\mu}_y, \Sigma_y)$$

where Σ_x and Σ_y are as in Equation (9). The normality is not essential for construction of this covariance structure.

The separable covariance structure brings complication. To account for this parsimonious structure in the generative model, we construct marginal factor matrices $\Sigma_{x,i}$ and $\Sigma_{y,j}$ that satisfy $\mathbf{V}_i^\top \Sigma_{x,i} \mathbf{V}_i = \mathbf{R}_x^{1/D_x} \mathbf{I}_{R_x}$ and $\mathbf{W}_j^\top \Sigma_{y,j} \mathbf{W}_j = \mathbf{R}_y^{1/D_y} \mathbf{I}_{R_y}$. Let $\mathbf{V}_i = \mathbf{Q}_{x,i} \mathbf{R}_{x,j}$ and $\mathbf{W}_j = \mathbf{Q}_{y,j} \mathbf{R}_{y,j}$ be the thin QR decompositions. Then

$$\Sigma_{x,i} = \mathbf{R}_x^{-D_x} \mathbf{Q}_{x,j} \mathbf{R}_{x,y}^{-T} \mathbf{R}_{x,j}^{-1} \mathbf{Q}_{x,i}^T + \mathbf{T}_{x,j} (\mathbf{I}_{D_i} - \mathbf{Q}_{x,i} \mathbf{Q}_{x,i}^T) \mathbf{T}_{x,j}^T$$

$$\Sigma_{y,j} = \mathbf{R}_y^{-D_y} \mathbf{Q}_{y,j} \mathbf{R}_{y,j}^{-T} \mathbf{R}_{y,j}^{-1} \mathbf{Q}_{y,j}^T + \mathbf{T}_{y,j} (\mathbf{I}_{q_j} - \mathbf{Q}_{y,j} \mathbf{Q}_{y,j}^T) \mathbf{T}_{y,j}^T$$

satisfy the conditions with arbitrary $\mathbf{T}_{x,i} \in \mathbb{R}^{p_i \times p_i}$ and $\mathbf{T}_{y,j} \in \mathbb{R}^{q_j \times q_j}$. And we have the desired property

$$\begin{aligned} & \text{vec}(\mathbf{V})^T (\Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1}) \text{vec}(\mathbf{V}) \\ &= \mathbf{1}_{R_x}^T (\mathbf{V}_{D_x} \circ \dots \circ \mathbf{V}_1)^T (\Sigma_{x,D_x} \otimes \dots \otimes \Sigma_{x,1}) (\mathbf{V}_{D_x} \circ \dots \circ \mathbf{V}_1) \mathbf{1}_{R_x} \\ &= \mathbf{1}_{R_x}^T (*, \mathbf{V}_i^T \Sigma_{x,i} \mathbf{V}_i) \mathbf{1}_x \\ &= \mathbf{R}_x^{-1} \mathbf{1}_{R_x}^T \mathbf{R}_x \mathbf{1}_{R_x} \\ &= 1. \end{aligned}$$

Similarly, $\text{vec}(\mathbf{W})^T (\Sigma_{y,D_y} \otimes \dots \otimes \Sigma_{y,1}) \text{vec}(\mathbf{W}) = 1$.

7 | NUMERICAL EXPERIMENTS

We use the generative model described in Section 6 to evaluate the methods discussed in this paper: classic CCA, TCCA, scTCCA, sparse TCCA, and sparse TCCA with the separable covariance. We assess these methods on their ability to recover the true latent population parameters (\mathcal{V} , \mathcal{W}) used to generate i.i.d. samples of tensor data pairs $(\mathcal{X}_n, \mathcal{Y}_n)$ for $n \in [1000]$. In all examples, the true \mathcal{V} is a vector of length 100 with six entries set to 1 and the rest to 0. We use three different latent $\mathcal{W} \in \mathbb{R}^{64 \times 64}$ shown in Figure 1: $\mathcal{W}_{\text{rectangle}}$, $\mathcal{W}_{\text{cross}}$ and $\mathcal{W}_{\text{butterfly}}$. White pixels indicate values of 1, and black pixels indicate values of 0. The rectangle and cross-population canonical tensors $\mathcal{W}_{\text{rectangle}}$ and $\mathcal{W}_{\text{cross}}$ are low rank; specifically, they are Rank 1 and Rank 2, respectively. The butterfly population canonical tensor $\mathcal{W}_{\text{butterfly}}$ is a high rank. At the first glance, these illustrative examples do not come across as challenging estimation problems in the high dimension-low sample size regime, but if we were to vectorize the data and perform CCA, the number of parameters to fit is $100 + 64^4 = 4196$ whereas the sample size is 1,000. Because the number of parameters exceeds the sample size, the sample covariance matrix is singular. Consequently, we add a small ridge term 10^{-3} to sample covariance matrices so that the generalized eigenvalue problem has a unique solution. The code for generating the simulation results is provided in the Supporting Information.

7.1 | Evaluation criteria and selection of tuning parameter

Canonical tensors can only be estimated up to a scaling factor. Thus, we measure estimation accuracy by the angle between population canonical tensors used to generate the data $(\mathcal{V}, \mathcal{W})$ and estimated canonical tensors $(\widehat{\mathcal{V}}, \widehat{\mathcal{W}})$, as follows:

$$\angle(\mathcal{V}, \widehat{\mathcal{V}}) = \frac{\langle \mathbf{v}, \widehat{\mathbf{v}} \rangle}{\|\mathbf{v}\|_2 \|\widehat{\mathbf{v}}\|_2}, \text{ and } \angle(\mathcal{W}, \widehat{\mathcal{W}}) = \frac{\langle \mathbf{w}, \widehat{\mathbf{w}} \rangle}{\|\mathbf{w}\|_2 \|\widehat{\mathbf{w}}\|_2}.$$

The angles can take on values from -1 to 1 where an angle closer to 1 indicates better recovery of the true canonical tensors.

We use k -fold cross-validation to perform model selection, for example, selection of the tensor rank or the sparsity level. We split our data into K equally sized groups; for $k \in [K]$, we estimate a pair $(\mathcal{V}_{-k}, \mathcal{W}_{-k})$ on all but the k th fold of data for a sequence of models of varying complexity, $\mathcal{M}_1, \mathcal{M}_2, \dots$, where model \mathcal{M}_1 has a fixed pair of tensor ranks R_x and R_y and sparsity inducing parameter λ . We denote the fitted canonical correlation using model \mathcal{M}_1 by $\widehat{\rho}_{-k}(\mathcal{V}_{-k}, \mathcal{W}_{-k}; \mathcal{M}_1)$. Using $(\mathcal{V}_{-k}, \mathcal{W}_{-k})$, we compute the canonical correlation on the held out j th fold, which we denote $\widehat{\rho}_{-j}(\mathcal{V}_{-k}, \mathcal{W}_{-k})$. We choose the model that minimizes the average discrepancy between the empirical canonical correlations on the training sets and testing sets (Waaijenborg, Verselewe de Witt Hamer, & Zwinderman, 2008).

$$\widehat{\mathcal{M}} = \mathcal{M}_{\text{argmin}} \frac{1}{K} \sum_{k=1}^K |\widehat{\rho}_k(\mathcal{V}_{-k}, \mathcal{W}_{-k}) - \widehat{\rho}_{-k}(\mathcal{V}_{-k}, \mathcal{W}_{-k}; \mathcal{M}_i)|.$$

We then use all the data to estimate $(\mathcal{V}, \mathcal{W})$ using model $\widehat{\mathcal{M}}$.

7.2 | Results

Figure 2 compares estimation accuracy of various methods over 100 replicates. First of all, we observe that the various versions of TCCA outperform the CCA on the vectorized data for all three choices of the latent canonical tensor \mathcal{W} . Second, we notice that there appears to be some overfitting in the case of rectangle and cross problems. In the rectangle problem, where the population canonical tensor $\mathcal{W}_{\text{rectangle}}$ is a Rank 1 tensor, all four TCCA methods show the best performance when we fix the rank R_y for estimating \mathcal{W} to be 1, and the performance deteriorates as higher ranks are used. The same trend can be seen in the result of the cross problem where the population canonical tensor $\mathcal{W}_{\text{cross}}$ is a Rank 2 tensor. All TCCA methods give smaller values of angle when $R_y = 3$ is used compared with when $R_y = 2$ is used. In contrast, in the case of the butterfly image, the rank of population canonical tensor $\mathcal{W}_{\text{butterfly}}$ is much greater than 3. Thus, we do not expect an overfitting problem as well as good estimation performance as in the low-rank cases. We confirm that in Figure 2, calculated angles from the butterfly problem are lower than those from rectangles or cross problem.

Another interesting observation is that results from TCCA methods with a separable covariance structure show better performance when we use a higher value for the rank parameter R_y than the true value than do other two models in cross and butterfly problems. This results implies that assuming a separable covariance structure improves the estimation accuracy in the two problems. This can be explained by the true image \mathcal{X} , which has a symmetric structure. Due to this special structure, we possibly expect improved performance from the more parsimonious model. However, this structured model does not have much effect on the rectangle problem. This may be because the true canonical tensor \mathcal{W} already possess relatively few parameters. Note that the rectangle image is also symmetric but Rank 1, which has very few parameters.

Table 1 shows computation times taken for each method. In most cases, the computation times of TCCA methods are better or comparable with those of the CCA method. Especially, there is huge improvement when the underlying true canonical tensor has a low rank. Also, sparse models take less time than do nonsparse models, and separable covariance structure models take less time than do the models without the assumption for computation in general, which is expected.

8 | DISCUSSION

We have proposed a TCCA approach for finding the relationship between linear combinations of two tensor datasets. Our method combines the classic CCA approach and the low-rank tensor decomposition to reduce the vast dimensionality of tensor parameters. The proposed estimation algorithm scales well with the tensor data size and is easy to implement using existing statistical software. Our algorithms also support convergence guarantees, properties arguably lacking in the alternatives in the literature. We briefly highlight some open problems for further investigation.

In our illustrative examples, we did not incorporate a procedure for selecting the rank parameter. A main motivation for the low-rank formulation of the canonical tensors is reduced computation. Thus, one may wish to choose as high a rank parameter as one's computational budget allows. Nonetheless, we leave for future work a more principled approach to rank selection. One promising angle of pursuit would be to leverage the equivalence between CCA and a least squares problem (Sun, Ji, & Ye, 2008) and then derive an information criterion, a strategy commonly used to select the rank parameter in several recently proposed tensor estimation procedures (Zhou et al., 2013; Sun et al., 2017; Sun & Li, 2019).

An additional direction for future work is to generalize our TCCA framework to handle the analysis of multiple tensor datasets. There are numerous proposed methods to extend CCA for multiple vector-measurement datasets (Carroll, 1968; Kettenring, 1971; Hanafi, 2007; Witten & Tibshirani, 2009; Luo et al., 2015) but not for multiple tensor datasets to the best of our knowledge.

Another problem not tackled in the paper includes verifying the consistency of the proposed estimation procedures. Our clear specification of the population models provides a

framework for studying the consistency property in both the large n fixed p (Zhou et al., 2013) and the large n diverging p settings (Zhang, Li, Zhou, Zhou, & Shen, 2019). Finally, we are currently investigating our methods on a real dataset.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

REFERENCES

- Bach FR, & Jordan MI (2006). A probabilistic interpretation of canonical correlation analysis, Technical Report.
- Bergstra J, & Bengio Y (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb), 281–305.
- Bickel PJ, & Levina E (2008). Covariance regularization by thresholding. *The Annals of Statistics*, 36(6), 2577–2604.
- Bickel PJ, & Levina E (2008). Regularized estimation of large covariance matrices. *The Annals of Statistics*, 36(1), 199–227.
- Cai TT, & Yuan M (2012). Adaptive covariance matrix estimation through block thresholding. *The Annals of Statistics*, 40(4), 2014–2042.
- Carroll JD (1968). Generalization of canonical correlation analysis to three or more sets of variables, *Proceedings of the 76th Annual Convention of the American Psychological Association*, Washington, DC: Vol. 3, pp. 227–228.
- Carroll JD, & Chang J-J (1970). Analysis of individual differences in multidimensional scaling via an N-way generalization of “Eckart-Young” decomposition. *Psychometrika*, 35, 283–319.
- Gang L, Yong Z, Yan-Lei L, & Jing D (2011). Three dimensional canonical correlation analysis and its application to facial expression recognition, *International Conference on Intelligent Computing and Information Science Berlin, Heidelberg: Springer*, pp. 56–61.
- González I, Déjean S, Martin P, & Baccini A (2008). CCA: An R package to extend canonical correlation analysis. *Journal of Statistical Software*, 23(1), 1–14.
- Hanafi M (2007). PLS path modelling: Computation of latent variables with the estimation mode b. *Computational Statistics*, 22(2), 275–292.
- Harshman RA (1970). Foundations of the PARAFAC procedure: Models and conditions for an “explanatory” multi-modal factor analysis, *UCLA Working Papers in Phonetics*, 16, 1–84.
- Hoff PD (2011). Separable covariance arrays via the Tucker product, with applications to multivariate relational data. *Bayesian Anal*, 6(2), 179–196.
- Hotelling H (1936). Relations between two sets of variates. *Biometrika*, 321–377.
- Kettenring JR (1971). Canonical analysis of several sets of variables. *Biometrika*, 58(3), 433–451.
- Kolda TG, & Bader BW (2009). Tensor decompositions and applications. *SIAM Review*, 51(3), 455–500.
- Kubokawa T, Srivastava MS, et al. (2013). Optimal ridge-type estimators of covariance matrix in high dimension: CIRJE, Faculty of Economics, University of Tokyo.
- Ledoit O, & Wolf M (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis*, 88(2), 365–411.
- Ledoit O, & Wolf M (2012). Nonlinear shrinkage estimation of large-dimensional covariance matrices. *The Annals of Statistics*, 40(2), 1024–1060.
- Lee SH, & Choi S (2007). Two-dimensional canonical correlation analysis. *Signal Processing Letters, IEEE*, 14(10), 735–738.
- Lu H (2013). Learning canonical correlations of paired tensor sets via tensor-to-vector projection, *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence, IJCAI '13 Beijing, China: AAAI Press*, pp. 1516–1522.

- Luo Y, Tao D, Ramamohanarao K, Xu C, & Wen Y (2015). Tensor canonical correlation analysis for multi-view dimension reduction. *IEEE Transactions on Knowledge and Data Engineering*, 27(11), 3111–3124.
- Ma Z (2013). Sparse principal component analysis and iterative thresholding. *The Annals of Statistics*, 41(2), 772–801.
- Srivastava MS, & Reid N (2012). Testing the structure of the covariance matrix with fewer observations than the dimension. *Journal of Multivariate Analysis*, 112, 156–171.
- Stein JL, Hua X, Lee S, Ho AJ, Leow AD, Toga AW & Thompson PM (2010). Voxelwise genome-wide association study (vGWAS). *Neuroimage*, 53(3), 1160–1174. [PubMed: 20171287]
- Sun L, Ji S, & Ye J (2008). A least squares formulation for canonical correlation analysis, Proceedings of the 25th International Conference on Machine Learning, ICML '08 New York, NY, USA: ACM, pp. 1024–1031.
- Sun WW, & Li L (2019). Dynamic tensor clustering. *Journal of the American Statistical Association*, in press.
- Sun WW, Lu J, Liu H, & Cheng G (2017). Provable sparse tensor decomposition. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3), 899–916.
- Tan KM, Wang Z, Liu H, & Zhang T (2018). Sparse generalized eigenvalue problem: Optimal statistical rates via truncated Rayleigh flow. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(5), 1057–1086.
- Vinod HD (1976). Canonical ridge and econometrics of joint production. *J Econ*, 4(2), 147–166.
- Waaijenborg S, Verselewe de Witt Hamer PC, & Zwinderman AH (2008). Quantifying the association between gene expressions and DNA-markers by penalized canonical correlation analysis. *Statistical Applications in Genetics and Molecular Biology*, 7(1), 3.
- Wang H (2010). Local two-dimensional canonical correlation analysis. *Signal Processing Letters, IEEE*, 17(11), 921–924.
- Wang Z, Gu Q, Ning Y, & Liu H (2015). High dimensional EM algorithm: Statistical optimization and asymptotic normality In Cortes C, Lawrence ND, Lee DD, Sugiyama M, & Garnett R (Eds.), *Advances in Neural Information Processing Systems 28*: Curran Associates, Inc, pp. 2521–2529.
- Wang S-J, Yan W-J, Sun T, Zhao G, & Fu X (2016). Sparse tensor canonical correlation analysis for micro-expression recognition. *Neurocomputing*, 214, 218–232.
- Witten DM, & Tibshirani RJ (2009). Extensions of sparse canonical correlation analysis with applications to genomic data. *Statistical Applications in Genetics and Molecular Biology*, 8(1), 1–27.
- Yan J, Zheng W, Zhou X, & Zhao Z (2012). Sparse 2-D canonical correlation analysis via low rank matrix approximation for feature extraction. *Signal Processing Letters, IEEE*, 19(1), 51–54.
- Zhang X, Li L, Zhou H, Zhou Y, & Shen D (2019). Tensor generalized estimating equations for longitudinal imaging analysis. *Statistica Sinica*, 29, 1977–2005. [PubMed: 32523321]
- 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. [PubMed: 24791032]

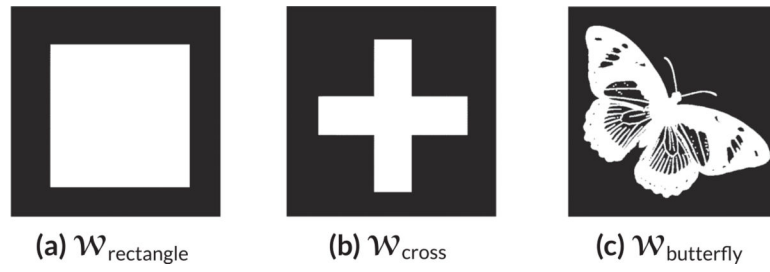


FIGURE 1.
True latent population \mathcal{W} canonical tensors for numerical experiments[dummy]

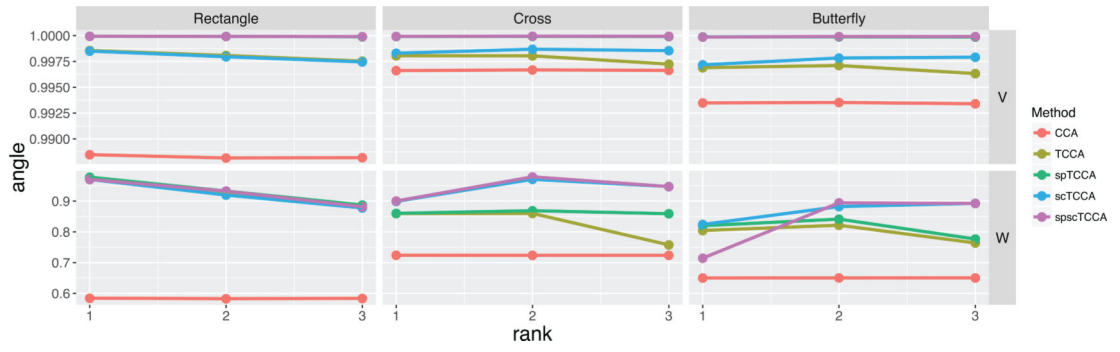


FIGURE 2. Angles between the recovered canonical vector/tensors and the true canonical vector/tensors. CCA, canonical correlation analysis; scTCCA, tensor canonical correlation analysis with separable covariance; spscTCCA, sparse tensor canonical correlation analysis with separable covariance; spTCCA, sparse tensor canonical correlation analysis; TCCA, tensor canonical correlation analysis

TABLE 1

Computation time: The mean run times are reported in seconds with the standard deviation in parentheses

Problem	Rank	CCA	TCCA	spTCCA	scTCCA	spscTCCA
Rectangle	1	22.90 (3.62)	7.52 (1.10)	7.36 (0.95)	3.14 (0.55)	3.09 (0.44)
	2	22.80 (3.18)	15.32 (2.51)	15.12 (2.17)	11.03 (2.03)	30.11 (8.78)
	3	21.85 (3.36)	24.07 (3.91)	18.61 (3.14)	15.89 (2.95)	37.09 (12.09)
Cross	1	18.49 (2.32)	8.05 (1.34)	7.78 (1.40)	3.39 (0.53)	13.19 (12.14)
	2	17.33 (2.43)	17.31 (3.32)	15.13 (5.02)	6.96 (1.11)	33.55 (12.78)
	3	17.10 (2.55)	24.44 (3.79)	54.86 (25.73)	10.04 (1.84)	9.63 (2.05)
Butterfly	1	18.92 (3.48)	7.08 (1.60)	61.18 (28.58)	2.90 (0.64)	31.66 (6.36)
	2	18.23 (2.92)	20.01 (4.53)	30.10 (14.87)	9.89 (2.15)	48.57 (13.14)
	3	19.99 (2.55)	32.27 (4.68)	26.09 (3.89)	11.84 (2.06)	13.58 (5.89)

Note. Experiments were performed on a computer cluster consisting of machines with Intel Xeon-based processors (I5 or I7 processors) with RAM ranging from 1,600 to 2,100 MHz. Abbreviations: CCA, canonical correlation analysis; scTCCA, tensor canonical correlation analysis with separable covariance; spscTCCA, sparse tensor canonical correlation analysis with separable covariance; spTCCA, sparse tensor canonical correlation analysis; TCCA, tensor canonical correlation analysis