Canonical Correlation Analysis of Datasets With a Common Source Graph

Jia Chen, Gang Wang ^(D), *Member, IEEE*, Yanning Shen ^(D), *Student Member, IEEE*, and Georgios B. Giannakis ^(D), *Fellow, IEEE*

Abstract—Canonical correlation analysis (CCA) is a powerful technique for discovering whether or not hidden sources are commonly present in two (or more) datasets. Its well-appreciated merits include dimensionality reduction, clustering, classification, feature selection, and data fusion. The standard CCA, however, does not exploit the geometry of the common sources, which may be available from the given data or can be deduced from (cross-) correlations. In this paper, this extra information provided by the common sources generating the data is encoded in a graph, and is invoked as a graph regularizer. This leads to a novel graph-regularized CCA approach, that is termed graph (g) CCA. The novel gCCA accounts for the graph-induced knowledge of common sources, while minimizing the distance between the wanted canonical variables. Tailored for diverse practical settings where the number of data is smaller than the data vector dimensions, the dual formulation of gCCA is developed too. One such setting includes kernels that are incorporated to account for nonlinear data dependencies. The resultant graph-kernel CCA is also obtained in closed form. Finally, corroborating image classification tests over several real datasets are presented to showcase the merits of the novel linear, dual, and kernel approaches relative to competing alternatives.

Index Terms—Dimensionality reduction, correlation analysis, signal processing over graphs, Laplacian regularization, generalized eigen-decomposition.

I. INTRODUCTION

I N MANY fields, exploratory data analysis depends critically on dimensionality reduction, a process to discover compact representations of large volumes of high-dimensional data [24]. Dimensionality reduction has been a crucial first step to obtain tractable learning tasks, such as classification, clustering, and regression [18], [24]. Principal component

J. Chen, Y. Shen, and G. B. Giannakis are with the Digital Technology Center and the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: chen5625@umn.edu; shenx513@umn.edu; georgios@umn.edu).

G. Wang is with the Digital Technology Center and the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN 55455 USA, and also with the State Key Lab of Intelligent Control and Decision of Complex Systems, Beijing Institute of Technology, Beijing 100081, China (e-mail: gangwang@umn.edu).

Color versions of one or more of the figures in this paper are available online at http://ieeexplore.ieee.org.

Digital Object Identifier 10.1109/TSP.2018.2853130

analysis (PCA) is arguably the most widely used dimensionality reduction method, finding low-dimensional representations from high-dimensional data points while preserving most of the data variance [19]. Nonetheless, ordinary PCA presumes that data vectors lie close to a hyperplane - a gross geometrical approximation for several datasets. Locally linear embedding on the other hand, preserves linear relationships between neighboring data [24], while Laplacian eigenmaps ensure that data close in the original manifold are mapped to close by locations in the low-dimensional space, therefore aiming to preserve local distances [4].

Nonetheless, such dimensionality reduction methods deal with one dataset at a time. They are challenged when it comes to analyzing two (or more) datasets jointly. Moreover, they require all data vectors to have the same dimension. Canonical correlation analysis (CCA) is a well-known method for extracting low-dimensional representations from two datasets that can have different dimensions, while maximizing their correlations [16]. Although recent PCA variants such as discriminative PCA can deal with two datasets at a time, their goal is to extract the most discriminative features from the data of interest relative to the other [10]. Formally, CCA aims at finding latent low-dimensional common structure from a paired dataset collected from different views of the same entities, also known as common sources. Each view contains high-dimensional representations of the sources in a certain feature space. For example, images of an individual captured by two cameras can be interpreted as two different views of this individual (here playing the role of a source). The ability of CCA to handle multiple datasets of different dimensions is a key enabler in diverse tasks such as multi-mode data fusion, where the need arises to fuse information from different domains [15]. Ever since its proposition [16], CCA benefits have been documented in diverse applications, such as blind source separation, brain imaging, clustering and classification, word embedding, and natural language processing, to name a few [12], [13], [15], [32].

To account for nonlinearities present in the data, kernel and deep CCA generalizations have also been developed based on kernels or deep neural networks [1], [15]. Sparse CCA looking for sparse canonical vectors was investigated by [33]. Multiview CCA on the other hand, generalizes ordinary CCA to handle data from more than two modalities [15]. Even though CCA solutions can be found via generalized eigen-decomposition, the resultant computational complexity may not scale well with

1053-587X © 2018 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.

Manuscript received March 27, 2018; revised June 24, 2018; accepted June 25, 2018. Date of publication July 9, 2018; date of current version July 18, 2018. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Olivier Lezoray. This work was partially supported by NSF under Grants 1500713 and 1514056. This paper was presented in part at the IEEE Statistical Signal Processing Workshop, Freiburg, Germany, June 10–13, 2018. (*Corresponding author: Gang Wang.*)

the problem dimensionality. This motivated decentralized CCA alternatives [9].

However, all aforementioned PCA and CCA tools do not exploit structural graph-induced information on the sources that may be available. Such information may be inferred from alternative views of the data, or it can be provided by the physics that dictates the underlying graph. Indeed, graph-aware dimensionality reduction methods have lately demonstrated promising performance [17], [25]–[27], [30].

Building on recent advances in graph-aware dimensionality reduction [17], [27], the present paper introduces a neat link between graph embedding and canonical correlations, by putting forward a novel graph (g) CCA approach. Our gCCA pursues maximally correlated linear projections, while also leveraging statistical dependencies due to the common sources hidden in the paired dataset. The underlying source graph encoding these dependencies can be either given, or be constructed based on prior knowledge. When the number of data samples is smaller than the data vector dimensions, we advocate the graph dual (gd) CCA. Relative to gCCA, our gdCCA not only bypasses the inversion of ill-conditioned data covariance matrices, but also incurs lower complexity in high-dimensional setups. To further account for nonlinearities, we also develop what we term graph kernel (gK) CCA. Interestingly, solutions to all three gCCA variants can be found analytically through generalized eigenvalue decompositions.

Different from [7], [35], where CCA was regularized by two graph Laplacians separately per view, gCCA here jointly leverages a single graph induced by the common sources. This is of major practical importance, e.g., in brain mapping, where besides functional magnetic resonance imaging (MRI) and diffusion-weighted MRI data collected at different brain regions, one has also access to the connectivity patterns among these regions [8]. Finally, numerical tests on several real-world datasets are presented to corroborate the merits of our proposed approaches for classification tasks over their competing alternatives.

The rest of this paper is structured as follows. Upon introducing the standard CCA in Section II, our gCCA is motivated, and derived in Section III. Its dual counterpart is developed in Section IV. Generalizing linear gCCA variants, the kernel version of gCCA is devised in Section V. Numerical tests on several real-world datasets are presented in Section VI, and the paper is concluded in Section VII.

Notation: Bold uppercase (lowercase) letters denote matrices (column vectors). Operators $\text{Tr}(\cdot)$, $(\cdot)^{-1}$ and $(\cdot)^{\top}$ are matrix trace, inverse and transpose, respectively; $\|\cdot\|_2$ stands for the ℓ_2 -norm of vectors; **0** is an all-zero vector whose dimension is clear from the context; $\langle \mathbf{a}, \mathbf{b} \rangle$ denotes the inner product of vectors **a** and **b**; and **I** represents the identity matrix of suitable size.

II. PRELIMINARIES

Consider two datasets $\{\mathbf{x}_i\}_{i=1}^N$ and $\{\mathbf{y}_i\}_{i=1}^N$ with corresponding dimensionality D_x and D_y , collected from two different views of the same sources $\{\mathbf{s}_i \in \mathbb{R}^{\rho}\}_{i=1}^N$ with possibly $\rho \ll \min\{D_x, D_y\}$. CCA amounts to finding low-dimensional subspaces $\mathbf{U} \in \mathbb{R}^{D_x \times d}$ and $\mathbf{V} \in \mathbb{R}^{D_y \times d}$ with $d \leq \rho$, such that the Euclidean distance between the low-dimensional representations $\{\mathbf{U}^{\top}\mathbf{x}_i\}$ and $\{\mathbf{V}^{\top}\mathbf{y}_i\}$ is minimized. Assume without loss of generality that both datasets are centered, meaning their corresponding sample means have been removed from the datasets. For ease of exposition, this section focuses on d = 1 first, while generalization to $d \geq 2$ will be discussed later. CCA solves the following problem

$$(\mathbf{u}^*, \mathbf{v}^*) := \arg\min_{\mathbf{u}, \mathbf{v}} \frac{1}{N} \sum_{i=1}^N \left(\mathbf{u}^\top \mathbf{x}_i - \mathbf{v}^\top \mathbf{y}_i \right)^2$$
 (1a)

where $\mathbf{u} \in \mathbb{R}^{D_x}$ and $\mathbf{v} \in \mathbb{R}^{D_y}$ are also termed a canonical pair. To ensure unique nonzero solutions however, the ensuing standard constraints are imposed

$$\mathbf{u}^{\top} \boldsymbol{\Sigma}_x \mathbf{u} = 1, \quad \text{and} \quad \mathbf{v}^{\top} \boldsymbol{\Sigma}_y \mathbf{v} = 1$$
 (1b)

where $\Sigma_x := (1/N) \sum_{i=1}^{N} \mathbf{x}_i \mathbf{x}_i^{\top}$ and $\Sigma_y := (1/N) \sum_{i=1}^{N} \mathbf{y}_i \mathbf{y}_i^{\top}$ denote the sample covariance matrices of $\{\mathbf{x}_i\}$ and $\{\mathbf{y}_i\}$, respectively. Projections $\{\mathbf{x}_i^{\top}\mathbf{u}^*\}_{i=1}^N$ and $\{\mathbf{y}_i^{\top}\mathbf{v}^*\}_{i=1}^N$ form a pair of canonical variables, which can be interpreted as lowdimensional approximations of the common sources $\{\mathbf{s}_i\}_{i=1}^N$.

After simple manipulations, (1) leads to the following popular formulation of CCA [15]

$$(\mathbf{u}^*, \mathbf{v}^*) := \arg \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}^\top \boldsymbol{\Sigma}_{xy} \mathbf{v}$$
 (2a)

s. to
$$\mathbf{u}^{\top} \boldsymbol{\Sigma}_x \mathbf{u} = 1$$
, and $\mathbf{v}^{\top} \boldsymbol{\Sigma}_y \mathbf{v} = 1$ (2b)

where $\Sigma_{xy} := (1/N) \sum_{i=1}^{N} \mathbf{x}_i \mathbf{y}_i^{\top}$ is the sample crosscovariance matrix of $\{\mathbf{x}_i\}$ and $\{\mathbf{y}_i\}$.

Using Lagrange duality theory, the solution of (2) will be given next in analytical form. To this end, letting $\lambda, \mu \in \mathbb{R}$ be the dual variables associated with the two constraints in (2b), one can write the Lagrangian as

$$\mathcal{L}(\mathbf{u}, \mathbf{v}; \lambda, \mu) = \mathbf{u}^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{v} - \lambda(\mathbf{u}^{\top} \boldsymbol{\Sigma}_{x} \mathbf{u} - 1) - \mu(\mathbf{v}^{\top} \boldsymbol{\Sigma}_{y} \mathbf{v} - 1).$$

At the optimum $(\mathbf{u}^*, \mathbf{v}^*)$, the KKT conditions assert that

$$\Sigma_{xy}\mathbf{v}^* = 2\lambda^*\Sigma_x\mathbf{u}^*, \qquad (\mathbf{u}^*)^\top\Sigma_x\mathbf{u}^* = 1$$
 (3a)

$$\boldsymbol{\Sigma}_{xy}^{\top} \mathbf{u}^* = 2\mu^* \boldsymbol{\Sigma}_y \mathbf{v}^*, \qquad (\mathbf{v}^*)^{\top} \boldsymbol{\Sigma}_y \mathbf{v}^* = 1.$$
(3b)

Left-multiplying the first equations in (3a) and (3b) by $(\mathbf{u}^*)^{\top}$ and $(\mathbf{v}^*)^{\top}$, respectively, leads to $(\mathbf{u}^*)^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{v}^* = 2\lambda^* = 2\mu^*$. Hence, solving (2) reduces to solving the generalized eigenvalue problem, see e.g., [15]

$$\begin{bmatrix} \boldsymbol{\Sigma}_{xy}^{\top} & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\Sigma}_{xy} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = 2\lambda \begin{bmatrix} \boldsymbol{0} & \boldsymbol{\Sigma}_{y} \\ \boldsymbol{\Sigma}_{x} & \boldsymbol{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}.$$
(4)

Maximizing the objective function (2a) is tantamount to finding the largest generalized eigenvalue $\lambda^* := \lambda_1$ in (4), and the optimal canonical vectors $[(\mathbf{u}^*)^\top (\mathbf{v}^*)^\top]^\top$ to (2) are obtained from the corresponding generalized eigenvector.

Consider the generalization of (2) to $d \leq \min(D_x, D_y)$ pairs of canonical vectors, say $\{(\mathbf{u}_i, \mathbf{v}_i)\}_{i=1}^d$. Upon letting $\mathbf{u}_1^* := \mathbf{u}^*$ and $\mathbf{v}_1^* := \mathbf{v}^*$, one can iteratively solve

$$\max_{\mathbf{u}_k, \mathbf{v}_k} \mathbf{u}_k^\top \boldsymbol{\Sigma}_{xy} \mathbf{v}_k \tag{5a}$$

s. to
$$\mathbf{u}_k^{\top} \boldsymbol{\Sigma}_x \mathbf{u}_k = 1$$
, $\mathbf{v}_k^{\top} \boldsymbol{\Sigma}_y \mathbf{v}_k = 1$ (5b)

$$\mathbf{u}_k^{\top} \boldsymbol{\Sigma}_x \mathbf{u}_i^* = 0, \quad \mathbf{v}_k^{\top} \boldsymbol{\Sigma}_y \mathbf{v}_i^* = 0 \tag{5c}$$

$$\forall i = 1, 2, \dots, k-1$$
 (5d)

for k = 2, 3, ..., d. For brevity, let us concatenate the d canonical vectors $\{\mathbf{u}_i\}$ and $\{\mathbf{v}_i\}$ to form matrices $\mathbf{U} \in \mathbb{R}^{D_x \times d}$ and $\mathbf{V} \in \mathbb{R}^{D_y \times d}$ accordingly, and rewrite (5) in the following compact form

$$\max_{\mathbf{U},\mathbf{V}} \operatorname{Tr}(\mathbf{U}^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{V})$$
(6a)

s. to
$$\mathbf{U}^{\top} \boldsymbol{\Sigma}_x \mathbf{U} = \mathbf{I}$$
, and $\mathbf{V}^{\top} \boldsymbol{\Sigma}_y \mathbf{V} = \mathbf{I}$ (6b)

which yields simultaneously multiple canonical vectors. As deduced earlier, the *m*-th columns of maximizers U^* and V^* of (6) correspond to the left and right generalized eigenvectors of (4) associated with the *m*-th largest generalized eigenvalue, respectively.

III. CCA OVER GRAPHS

In diverse applications, the common sources $\{s_i\}_{i=1}^N$ may be viewed as nodal vectors of a graph having N nodes. This structural prior information can be leveraged when finding the canonical vectors. In this paper, this extra knowledge of common sources is encoded in a graph, and will be embodied in the canonical variables through graph regularization.

We outline some basics of the graph theory first. A graph is represented by a tuple $\mathcal{G} = \{\mathcal{N}, \mathcal{W}\}$, where $\mathcal{N} := \{1, 2, ..., N\}$ is the vertex set, and $\mathcal{W} := \{w_{ij}\}_{(i,j)\in\mathcal{N}\times\mathcal{N}}$ stacks up edge weights w_{ij} over all vertex pairs (i, j). For ease of exposition, this paper focuses on undirected graphs, for which $w_{ij} = w_{ji}$ for all $i, j \in \mathcal{N}$. Moreover, a graph is said to be unweighted if all w_{ij} 's take binary values 0 or 1. Upon forming the so-called weighted adjacency matrix $\mathbf{W} \in \mathbb{R}^{N\times N}$ with its (i, j)-th entry being w_{ij} , and defining $d_i := \sum_{j=1}^N w_{ij}$, the Laplacian matrix of graph \mathcal{G} is given by

$$\mathbf{L}_{G} := \mathbf{D} - \mathbf{W} \in \mathbb{R}^{N \times N} \tag{7}$$

where the diagonal matrix $\mathbf{D} \in \mathbb{R}^{N \times N}$ holds ordered entries $\{d_i\}_{i=1}^N$ on its diagonal.

Having introduced basic graph notation, we present a neat link between canonical correlations and graph embedding next. Consider for instance a graph \mathcal{G} with adjacency matrix \mathbf{W} , over which the underlying sources $\{\mathbf{s}_i\}_{i=1}^N$ are assumed to be smooth. In other words, vectors $(\mathbf{s}_i, \mathbf{s}_j)$ residing on two connected nodes $i, j \in \mathcal{G}$ are deemed close to each other in Euclidean distance. As remarked earlier, canonical variables $\mathbf{u}^\top \mathbf{x}_i$ and $\mathbf{v}^\top \mathbf{y}_j$ are accordingly one-dimensional approximates of \mathbf{s}_i and \mathbf{s}_j . Building on this fact, let us now focus on the weighted sum of distances between any two pairs of canonical variables from $\{\mathbf{u}^\top \mathbf{x}_i\}_{i=1}^N$ and $\{\mathbf{v}^\top \mathbf{y}_i\}_{i=1}^N$ over \mathcal{G} , namely the quadratic term

$$\sum_{i=1}^{N} \sum_{j=1}^{N} w_{ij} \left(\mathbf{u}^{\top} \mathbf{x}_{i} - \mathbf{v}^{\top} \mathbf{y}_{j} \right)^{2}.$$
(8)

It is clear that by minimizing (8) over \mathbf{u} and \mathbf{v} , canonical variables $\mathbf{u}^{\top}\mathbf{x}_i$ and $\mathbf{v}^{\top}\mathbf{y}_j$ corresponding to adjacent nodes $i, j \in \mathcal{G}$ with large edge weights w_{ij} will be promoted to stay close to each other. As such, invoking this term as a regularizer accounts for the additional graph knowledge of the common sources, while maximizing the linear correlation coefficient between the canonical variables, yielding

$$\min_{\mathbf{u},\mathbf{v}} \quad \frac{1}{2N} \sum_{i=1}^{N} \left(\mathbf{u}^{\mathsf{T}} \mathbf{x}_{i} - \mathbf{v}^{\mathsf{T}} \mathbf{y}_{i} \right)^{2} + \frac{\gamma}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} w_{ij} \left(\mathbf{u}^{\mathsf{T}} \mathbf{x}_{i} - \mathbf{v}^{\mathsf{T}} \mathbf{y}_{j} \right)^{2}$$

s. to $\mathbf{u}^{\mathsf{T}} \mathbf{\Sigma}_{x} \mathbf{u} = 1$, and $\mathbf{v}^{\mathsf{T}} \mathbf{\Sigma}_{y} \mathbf{v} = 1$

in which $\gamma \ge 0$ is a hyper-parameter that balances the distance between canonical variable estimates with their smoothness over \mathcal{G} . After expanding the squares and removing the constant terms, the problem at hand can be equivalently rewritten as

$$\max_{\mathbf{u},\mathbf{v}} \mathbf{u}^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{v} - \gamma \mathbf{u}^{\top} \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top} \mathbf{v} - \frac{\gamma}{2} \sum_{i=1}^{N} d_{i} \left(\mathbf{u}^{\top} \mathbf{x}_{i} - \mathbf{v}^{\top} \mathbf{y}_{i} \right)^{2}$$
(9a)

s. to $\mathbf{u}^{\top} \boldsymbol{\Sigma}_x \mathbf{u} = 1$, and $\mathbf{v}^{\top} \boldsymbol{\Sigma}_y \mathbf{v} = 1$ (9b)

where $\mathbf{X} := [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D_x \times N}$, and $\mathbf{Y} := [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N] \in \mathbb{R}^{D_y \times N}$.

Evidently, (9) is non-convex and is not amenable to efficient solvers due to the bilinear terms as well as the quadratic equality constraints. Even though block coordinate descent-type solvers can be employed, only convergence to a stationary point can be guaranteed in general [9]. Instead of coping with the objective function (9a) directly, we shall pursue a lower bound of it, which will turn out to afford an analytical solution.

Toward that end, it is easy to verify that with all $\{d_i \ge 0\}_{i=1}^N$, the following holds for all $\mathbf{u} \in \mathbb{R}^{D_x}$ and $\mathbf{v} \in \mathbb{R}^{D_y}$

$$\sum_{i=1}^{N} d_i \left(\mathbf{u}^{\top} \mathbf{x}_i - \mathbf{v}^{\top} \mathbf{y}_i \right)^2 \le 2 d_{\max} N \left(\mathbf{u}^{\top} \boldsymbol{\Sigma}_x \mathbf{u} + \mathbf{v}^{\top} \boldsymbol{\Sigma}_y \mathbf{v} \right) \quad (10)$$

where $d_{\max} := \max_{1 \le i \le N} d_i$, and the equality is achieved when $d_i = d_{\max}$ and $\mathbf{u}^\top \mathbf{x}_i = -\mathbf{v}^\top \mathbf{y}_i$ for all i = 1, 2, ..., N. Subsequently, we replace the last term in (9a) with the right-hand-side term, which contributes to a valid lower bound of (9a). Formally stated, we have the following reformulation.

Proposition 1: Replacing the sum in (9a) with its upper bound in (10) leads to an objective that lower bounds (9a). Merging and ignoring the constant terms due to the equality constraints (9b) leads to our novel gCCA formulation

$$\max_{\mathbf{u},\mathbf{v}} \ \mathbf{u}^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{v} - \gamma \mathbf{u}^{\top} \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top} \mathbf{v}$$
(11a)

s. to
$$\mathbf{u}^{\top} \boldsymbol{\Sigma}_x \mathbf{u} = 1$$
, and $\mathbf{v}^{\top} \boldsymbol{\Sigma}_y \mathbf{v} = 1$. (11b)

Clearly, when $\gamma = 0$, our gCCA finds (\mathbf{u}, \mathbf{v}) that only maximizes the linear correlation between the pair of canonical variables. In this case, no graph knowledge is exploited, and our gCCA reduces to the standard CCA. With γ increasing gradually, gCCA accounts progressively for extra graph information of the common sources when finding the canonical variables. Note that quantifying the gap between (9) and (11) is

- Input: {x_i}^N_{i=1}, {y_i}^N_{i=1}, d, W, and γ.
 Form (cross-)covariance matrices, Σ_x, Σ_y and Σ_{xy}.
- 3: **Build** $L_{\mathcal{G}}$ using (7).
- 4: Perform SVD on $\Sigma_x^{-1/2} \left(\Sigma_{xy} \gamma \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top} \right) \Sigma_y^{-1/2}$
- 5: Extract the first d leading left and right singular vectors to obtain $\overline{\mathbf{U}}^*$ and $\overline{\mathbf{V}}^*$, respectively.
- 6: Compute $\mathbf{U}^* = \boldsymbol{\Sigma}_x^{-1/2} \overline{\mathbf{U}}^*$ and $\mathbf{V}^* = \boldsymbol{\Sigma}_y^{-1/2} \overline{\mathbf{V}}^*$.
- 7: Output: U^* and V^* .

challenging, as it depends on the collected data realizations $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$.

Next, let us consider multiple canonical pairs $\{(\mathbf{u}_i, \mathbf{v}_i)\}_{i=1}^d$, and collect them to form matrices $\mathbf{U} := [\mathbf{u}_1 \cdots \mathbf{u}_d]$ and $\mathbf{V} :=$ $[\mathbf{v}_1 \cdots \mathbf{v}_d]$. We can then generalize gCCA in (11) to $d \ge 2$ as

$$\max_{\mathbf{U},\mathbf{V}} \operatorname{Tr} \left(\mathbf{U}^{\top} \boldsymbol{\Sigma}_{xy} \mathbf{V} - \gamma \mathbf{U}^{\top} \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top} \mathbf{V} \right)$$
(12a)

s. to
$$\mathbf{U}^{\top} \boldsymbol{\Sigma}_x \mathbf{U} = \mathbf{I}$$
, and $\mathbf{V}^{\top} \boldsymbol{\Sigma}_y \mathbf{V} = \mathbf{I}$. (12b)

Interestingly, even with the extra graph-inducing regularization term, our gCCA in (12) still admits an analytical solution, under the standard assumption that data covariance matrices Σ_x and Σ_{y} are both nonsingular. For concreteness, the solution is summarized in the following result, and for self-contained presentation, its proof is provided in Appendix A.

Theorem 1: Given zero-mean data $\{\mathbf{x}_i \in \mathbb{R}^{D_x}\}_{i=1}^N$ and $\{\mathbf{y}_i \in \mathbb{R}^{D_y}\}_{i=1}^N$, suppose that $\mathbf{\Sigma}_x = (1/N)\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top$ and $\mathbf{\Sigma}_y = (1/N)\sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i^\top$ are nonsingular. Then the optimal solution $(\mathbf{U}^* \in \mathbb{R}^{D_x \times d}, \mathbf{V}^* \in \mathbb{R}^{D_y \times d})$ to the gCCA problem (12) with $d \leq \min(D_x, D_y)$, is given by

$$\mathbf{U}^* := \boldsymbol{\Sigma}_x^{-1/2} \bar{\mathbf{U}}^*, \quad \text{and} \quad \mathbf{V}^* := \boldsymbol{\Sigma}_y^{-1/2} \bar{\mathbf{V}}^* \qquad (13)$$

where the columns of $\bar{\mathbf{U}}^* \in \mathbb{R}^{D_x \times d}$ and $\bar{\mathbf{V}}^* \in \mathbb{R}^{D_y \times d}$ are the d left and right singular vectors of $\Sigma_x^{-1/2}(\Sigma_{xy})$ $\gamma \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top}) \boldsymbol{\Sigma}_{y}^{-1/2}$ associated with its *d* largest singular values. Moreover, the maximum objective value of (12a) is the sum of the d largest singular values.

Our proposed gCCA scheme is summarized in Alg. 1. A couple of remarks are now in order.

Remark 1: Different from our single regularizer in (12), the approaches in [7], [35] rely on two regularizers or two constraints involving graph priors $\mathbf{U}^{\top}\mathbf{X}\mathbf{L}_{\mathcal{G}_{w}}\mathbf{X}^{\top}\mathbf{U}$ and $\mathbf{V}^{\top}\mathbf{Y}\mathbf{L}_{\mathcal{G}_{u}}\mathbf{Y}^{\top}\mathbf{V}$ for the two-view data **X** and **Y**, respectively. However, the problem formulation in [35] does not admit an analytical solution. Although iterative algorithms can be used to solve the involved nonconvex optimization problem, only convergence to a stationary point can be ensured in general [6]. When the two datasets lie in two distinct graphs \mathcal{G}_x and \mathcal{G}_y , using the graph-Laplacian regularized constraints can improve standard CCA performance [5]. This approach is mainly suggested for semi-supervised learning, where Σ_{xy} is fully available. In contrast, (12) leverages the graph induced by the common sources, and our source graph regularizer $\mathbf{U}^{\top}\mathbf{X}\mathbf{L}_{G}\mathbf{Y}^{\top}\mathbf{V}$ directly exploits correlations between the low-dimensional approximations of common sources over \mathcal{G} . This is critical in certain practical setups, in which one has prior knowledge about the common sources besides the given datasets. In general, the graph of inter-dependent sources can be *a priori* provided by an 'expert' or dictated by the underlying physics, or, it can be learned from alternate views of the data [14], [29]. For example, in electric power networks, besides the power quantities observed, one has also access to the grid topology [20] capturing the connectivities between buses (substations) through power lines. Likewise, in brain networks, in addition to the functional MRI and diffusion-weighted MRI data collected at different brain regions, the connectivity patterns among these regions may also be available through other means [8]; see e.g., the UCLA Multimodal Connectivity Database¹, a repository for researchers to publicly share the connectivity matrices derived from their research. Last but not the least, our proposed gCCA framework features an analytical solution.

Remark 2: To induce different graph properties, rather than relying on $\mathbf{L}_{\mathcal{G}}$, a family of graph regularizations of the form $r(\mathbf{L}_{\mathcal{G}}) := \sum_{i=1}^{N} r(\lambda_i^l) \boldsymbol{\nu}_i \boldsymbol{\nu}_i^{\top}$ can be also employed [31], where $r(\cdot): \mathbb{R} \to \mathbb{R}^+$ is a scalar function, and appropriate choices of $r(\lambda_i^l)$ are helpful for inducing diverse graph properties; while $\boldsymbol{\nu}_i \in \mathbb{R}^N$ is the eigenvector of $\mathbf{L}_{\mathcal{G}}$ associated with its *i*-th largest eigenvalue λ_i^l .

Remark 3: The hyper-parameter γ can be determined through e.g., the following two ways: i) cross-validation for supervised tasks, where one is given labeled training samples, and γ is chosen as the value optimizing empirical performance on the training samples; and, ii) a spectral clustering based approach that automatically selects the best γ value from a given set of candidate values, as in [10].

Remark 4: Bayesian approaches can also be adopted in capturing prior knowledge of the common sources (which in fact can be generic distributions). But when this distribution is not given a priori, which is true in most data analytics applications, estimating them requires typically a prohibitively large number of training samples (a.k.a. "curse of dimensionality"). In contrast, the pursued graph-regularized approach can capture much more subtle, intricate, and (even) non-metric inter-dependencies of the common sources. Consider for instance a binary adjacency matrix \mathbf{W} , and the matrix $\mathbf{W}\mathbf{W}^{\top}$ which reveals the number of common neighbors of any two nodes. The latter can be viewed as a covariance matrix of sources with a Gaussian prior distribution in a Bayesian setup, and can be obtained readily even without source realizations. When the graph is unknown, identifying an (even approximate) graph is often easier than estimating a distribution from limited number of data samples.

IV. DUAL CCA OVER GRAPHS

Similar to dual PCA [27], various practical scenarios involving high-dimensional data vectors, have $N \ll \min\{D_x, D_y\}$, in which case Σ_x and Σ_y become singular, and the results in Theorem 1 do not apply. Even though this rank deficiency can be remedied with appropriate Tikhonov regularization [15], the

¹http://umcd.humanconnectomeproject.org.

resultant computational complexity can be considerably higher than the alternative of investigating gCCA in the dual domain. In this direction, consider first expressing $\mathbf{u} \in \mathbb{R}^{D_x}$ and $\mathbf{v} \in \mathbb{R}^{D_y}$ in terms of their corresponding parts of the data matrices **X** and **Y** as

$$\mathbf{u} := \mathbf{X}\boldsymbol{\alpha}, \quad \text{and} \quad \mathbf{v} := \mathbf{Y}\boldsymbol{\beta}$$
 (14)

where $\alpha \in \mathbb{R}^N$ and $\beta \in \mathbb{R}^N$ are the so-termed dual vectors. Substituting (14) into (11) gives rise to our graph dual (gd) CCA formulation for one pair of canonical vectors

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta}} \boldsymbol{\alpha}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{Y}^{\top} \mathbf{Y} \boldsymbol{\beta} - \gamma \boldsymbol{\alpha}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{L}_{\boldsymbol{\beta}} \mathbf{Y}^{\top} \mathbf{Y} \boldsymbol{\beta}$$
(15a)

s. to
$$\boldsymbol{\alpha}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\alpha} = 1$$
 (15b)

$$\boldsymbol{\beta}^{\top} \mathbf{Y}^{\top} \mathbf{Y} \mathbf{Y}^{\top} \mathbf{Y} \boldsymbol{\beta} = 1.$$
(15c)

Similar to Section II, introducing variables $\lambda_x \in \mathbb{R}$ and $\lambda_y \in \mathbb{R}$ to be the Lagrange multipliers corresponding to constraints (15b) and (15c), respectively, one can write the Lagrangian for (15) as

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha},\,\boldsymbol{\beta};\,\lambda_x,\,\lambda_y) &:= -\boldsymbol{\alpha}^\top \mathbf{X}^\top \mathbf{X} (\mathbf{I} - \gamma \mathbf{L}_{\boldsymbol{\mathcal{G}}}) \mathbf{Y}^\top \mathbf{Y} \boldsymbol{\beta} \\ &+ \frac{\lambda_x}{2} (\boldsymbol{\alpha}^\top \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \boldsymbol{\alpha} - 1) + \frac{\lambda_y}{2} (\boldsymbol{\beta}^\top \mathbf{Y}^\top \mathbf{Y} \mathbf{Y}^\top \mathbf{Y} \boldsymbol{\beta} - 1). \end{aligned}$$

Setting derivatives of the Lagrangian with respect to α and β to zero further leads to

$$-\mathbf{X}^{\top}\mathbf{X}(\mathbf{I}-\gamma\mathbf{L}_{\mathcal{G}})\mathbf{Y}^{\top}\mathbf{Y}\boldsymbol{\beta}+\lambda_{x}\mathbf{X}^{\top}\mathbf{X}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\alpha}=\mathbf{0}$$
(16a)

$$-\mathbf{Y}^{\top}\mathbf{Y}(\mathbf{I}-\gamma\mathbf{L}_{\mathcal{G}})\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\alpha}+\lambda_{y}\mathbf{Y}^{\top}\mathbf{Y}\mathbf{Y}^{\top}\mathbf{Y}\boldsymbol{\beta}=\mathbf{0}.$$
 (16b)

Left-multiplying (16a) and (16b) by α^{\top} and β^{\top} , respectively, and subsequently subtracting the latter from the former, we arrive at

$$\lambda_x \boldsymbol{\alpha}^\top \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \boldsymbol{\alpha} - \lambda_y \boldsymbol{\beta}^\top \mathbf{Y}^\top \mathbf{Y} \mathbf{Y}^\top \mathbf{Y} \boldsymbol{\beta} = 0.$$
(17)

Taking into account (17), (15b), and (15c), it follows that at the optimal solution, we have $\lambda^* := \lambda_x^* = \lambda_y^*$. Supposing for now that $\mathbf{X}^\top \mathbf{X}$ and $\mathbf{Y}^\top \mathbf{Y}$ are nonsingular, we find

$$\boldsymbol{\alpha}^* := \frac{1}{\lambda^*} \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \left(\mathbf{Y}^\top \mathbf{Y} - \gamma \mathbf{L}_{\mathcal{G}} \mathbf{Y}^\top \mathbf{Y} \right) \boldsymbol{\beta}^*.$$
(18)

Plugging (18) into (16b) yields

$$\left(\mathbf{Y}^{\top}\mathbf{Y}\right)^{-1}\left(\mathbf{I}-\gamma\mathbf{L}_{\mathcal{G}}\right)^{2}\mathbf{Y}^{\top}\mathbf{Y}\boldsymbol{\beta}^{*}=(\lambda^{*})^{2}\boldsymbol{\beta}^{*}$$
 (19)

and similarly, one obtains that

$$\left(\mathbf{X}^{\top}\mathbf{X}\right)^{-1}\left(\mathbf{I}-\gamma\mathbf{L}_{\mathcal{G}}\right)^{2}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\alpha}^{*}=(\lambda^{*})^{2}\boldsymbol{\alpha}^{*}.$$
 (20)

The last two equalities show that α^* depends solely on **X**, and β^* solely on **Y**. This holds without any assumption about the paired dataset **X** and **Y** whatsoever. Furthermore, when $\gamma = 0$, both (19) and (20) lead to trivial solutions. However, recall that our goal is to extract relations between data **X** and **Y**. As with the dual CCA [15], in order to avoid such trivial solutions, we invoke two Tikhonov regularization terms that lead to our graph

dual (gd) CCA formulation

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta}} \ \boldsymbol{\alpha}^{\top} \left(\mathbf{X}^{\top} \mathbf{X} \mathbf{Y}^{\top} \mathbf{Y} - \gamma \mathbf{X}^{\top} \mathbf{X} \mathbf{L}_{\boldsymbol{\beta}} \mathbf{Y}^{\top} \mathbf{Y} \right) \boldsymbol{\beta}$$
(21a)

s. to
$$\boldsymbol{\alpha}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\alpha} + \epsilon \boldsymbol{\alpha}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\alpha} = 1$$
 (21b)

$$\boldsymbol{\beta}^{\top} \mathbf{Y}^{\top} \mathbf{Y} \mathbf{Y}^{\top} \mathbf{Y} \boldsymbol{\beta} + \epsilon \boldsymbol{\beta}^{\top} \mathbf{Y}^{\top} \mathbf{Y} \boldsymbol{\beta} = 1.$$
(21c)

Here, the coefficient $\epsilon > 0$ is a pre-selected penalty parameter. Appealing to Lagrange duality theory again, one arrives at

$$(\mathbf{I} - \gamma \mathbf{L}_{\mathcal{G}})\mathbf{Y}^{\top}\mathbf{Y}(\mathbf{Y}^{\top}\mathbf{Y} + \epsilon \mathbf{I})^{-1}(\mathbf{I} - \gamma \mathbf{L}_{\mathcal{G}})\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\alpha}^{*}$$

$$= (\lambda^{*})^{2}(\mathbf{X}^{\top}\mathbf{X} + \epsilon \mathbf{I})\boldsymbol{\alpha}^{*} \qquad (22a)$$

$$(\mathbf{I} - \gamma \mathbf{L}_{\mathcal{G}})\mathbf{X}^{\top}\mathbf{X}(\mathbf{X}^{\top}\mathbf{X} + \epsilon \mathbf{I})^{-1}(\mathbf{I} - \gamma \mathbf{L}_{\mathcal{G}})\mathbf{Y}^{\top}\mathbf{Y}\boldsymbol{\beta}^{*}$$

$$= (\lambda^{*})^{2}(\mathbf{Y}^{\top}\mathbf{Y} + \epsilon \mathbf{I})\boldsymbol{\beta}^{*} \qquad (22b)$$

suggesting that the maximizers α^* and β^* are accordingly the eigenvectors of (22a) and (22b) associated with the largest generalized eigenvalue $(\lambda_1^*)^2$. Moreover, the optimal objective function value in (21a) coincides with λ_1^* .

When looking for *d* pairs of dual vectors $\{(\alpha_i, \beta_i)\}_{i=1}^d$, which are collected to form matrices $\mathbf{A} := [\alpha_1 \cdots \alpha_d] \in \mathbb{R}^{N \times d}$ and $\mathbf{B} := [\beta_1 \cdots \beta_d] \in \mathbb{R}^{N \times d}$, our gdCCA becomes

$$\max_{\mathbf{A},\mathbf{B}} \operatorname{Tr} \left(\mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{Y}^{\top} \mathbf{Y} \mathbf{B} - \gamma \mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{L}_{\mathcal{G}} \mathbf{Y}^{\top} \mathbf{Y} \mathbf{B} \right)$$
(23a)

s. to
$$\mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{X}^{\top} \mathbf{X} \mathbf{A} + \epsilon \mathbf{A}^{\top} \mathbf{X}^{\top} \mathbf{X} \mathbf{A} = \mathbf{I}$$
 (23b)

$$\mathbf{B}^{\mathsf{T}}\mathbf{Y}^{\mathsf{T}}\mathbf{Y}\mathbf{Y}^{\mathsf{T}}\mathbf{Y}\mathbf{B} + \epsilon\mathbf{B}^{\mathsf{T}}\mathbf{Y}^{\mathsf{T}}\mathbf{Y}\mathbf{B} = \mathbf{I}$$
(23c)

for which the *i*-th columns of its optimal solutions \mathbf{A}^* and \mathbf{B}^* are accordingly provided by the generalized eigenvectors in (22a) and (22b) associated with the *i*-th largest generalized eigenvalue. Once \mathbf{A}^* , \mathbf{B}^* are found, the optimal canonical vectors sought can be obtained via (14) as $\mathbf{U}^* = \mathbf{X}\mathbf{A}^*$ and $\mathbf{V}^* = \mathbf{Y}\mathbf{B}^*$.

V. KCCA OVER GRAPHS

Although linear models are attractive due to their simplicity, they cannot capture complex nonlinear data dependencies that are common in real-world applications, including genomics [34], functional MRI [7], and acoustic feature learning [1].

Going beyond linearity, we generalize our linear models of CCA over graphs in Sections III and IV to take into account nonlinear relationships between data X and Y using kernel methods. In this context, a graph kernel (gK) CCA framework is developed. We begin with transforming the two datasets using two nonlinear functions to higher (possibly infinite) dimensional feature spaces, and subsequently find low-dimensional canonical variables. Specifically, let ϕ_x be a mapping from space \mathbb{R}^{D_x} to space \mathbb{R}^{D_h} (possibly with $D_h = \infty$). It is clear from (23) that both the objective and the constraints depend on the data X only through the similarities $\{\langle \mathbf{x}_i, \mathbf{x}_j \rangle\}_{i,j=1}^N$. Therefore, upon 'lifting' all data vectors $\{\mathbf{x}_i\}_{i=1}^N$ to obtain $\{\phi(\mathbf{x}_i)\}_{i=1}^N$, all similarities $\{\langle \mathbf{x}_i, \mathbf{x}_j \rangle\}_{i,j=1}^N$ can be readily replaced with $\{\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle\}_{i,j=1}^N$ can be computationally intractable due to the high-dimensionality.

To circumvent the cost of explicitly working in the highdimensional space, the so-called 'kernel trick' is employed [2]. To this end, we select some kernel function κ_x , such that $\kappa_x(\mathbf{x}_i, \mathbf{x}_j) := \langle \phi_x(\mathbf{x}_i), \phi_x(\mathbf{x}_j) \rangle$ for all i, j = 1, 2, ..., N, which form the (i, j)-th entries of the so-termed kernel matrix $\mathbf{\bar{K}}_x \in \mathbb{R}^{N \times N}$. Similarly, we can build the kernel matrix $\mathbf{\bar{K}}_y \in \mathbb{R}^{N \times N}$ for data \mathbf{Y} using a different kernel function κ_y . As in linear gCCA and gdCCA discussed is Sections III and IV, we require that the data in the mapped feature spaces $\{\phi_x(\mathbf{x}_i)\}_{i=1}^N$ and $\{\phi_y(\mathbf{y}_i)\}_{i=1}^N$ be centered, where $\phi_y(\mathbf{y}_i)$ is the nonlinear mapping for 'lifting' data \mathbf{y}_i to render kernel matrix \mathbf{K}_y . Using the kernel trick again, the required centering in the highdimensional space can be realized by centering the kernel matrix for data \mathbf{X} as

$$\mathbf{K}_{x}(i,j) := \bar{\mathbf{K}}_{x}(i,j) - \frac{1}{N} \sum_{\ell=1}^{N} \bar{\mathbf{K}}_{x}(\ell,j) - \frac{1}{N} \sum_{\ell=1}^{N} \bar{\mathbf{K}}_{x}(i,\ell) + \frac{1}{N^{2}} \sum_{m=1}^{N} \sum_{n=1}^{N} \bar{\mathbf{K}}_{x}(m,n)$$
(24)

and likewise for centering \mathbf{K}_y .

Upon replacing $\mathbf{X}^{\top}\mathbf{X}$ and $\mathbf{Y}^{\top}\mathbf{Y}$ in (23) with \mathbf{K}_x and \mathbf{K}_y , we arrive at our gKCCA

$$\max_{\mathbf{A},\mathbf{B}} \operatorname{Tr}(\mathbf{A}^{\top}\mathbf{K}_{x}\mathbf{K}_{y}\mathbf{B} - \gamma\mathbf{A}^{\top}\mathbf{K}_{x}\mathbf{L}_{\mathcal{G}}\mathbf{K}_{y}\mathbf{B})$$
(25a)

s. to
$$\mathbf{A}^{\top}\mathbf{K}_{x}^{2}\mathbf{A} + \epsilon\mathbf{A}^{\top}\mathbf{K}_{x}\mathbf{A} = \mathbf{I}$$
 (25b)

$$\mathbf{B}^{\top}\mathbf{K}_{y}^{2}\mathbf{B} + \epsilon \mathbf{A}^{\top}\mathbf{K}_{y}\mathbf{B} = \mathbf{I}.$$
 (25c)

It is clear that with properly selected kernel matrices K_x and K_y , gKCCA is able to capture nonlinear correlations between X and Y, while also leveraging the graph prior information of the common sources. Following the steps used to solve the gCCA problem (12), the solution to (25) is summarized in Theorem 2, with its proof deferred to Appendix B. The main steps of the gKCCA are listed in Alg. 2.

Theorem 2: If \mathbf{K}_x and \mathbf{K}_y are nonsingular, the optimal solutions \mathbf{A}^* and \mathbf{B}^* to (25) are given by

$$\mathbf{A}^* := \mathbf{K}_x^{-1/2} (\mathbf{K}_x + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{A}}^*$$
(26a)

$$\mathbf{B}^* := \mathbf{K}_y^{-1/2} (\mathbf{K}_y + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{B}}^*$$
(26b)

where matrices $\bar{\mathbf{A}}^* \in \mathbb{R}^{N \times d}$ and $\bar{\mathbf{B}}^* \in \mathbb{R}^{N \times d}$ collect as columns the top d left and right singular vectors of

$$\mathbf{C} := (\mathbf{K}_x + \epsilon \mathbf{I})^{-1/2} \mathbf{K}_x^{1/2} (\mathbf{I} - \gamma \mathbf{L}_{\mathcal{G}}) \mathbf{K}_y^{1/2} (\mathbf{K}_y + \epsilon \mathbf{I})^{-1/2}.$$
(26c)

Furthermore, the optimal objective value (25a) is the sum of the d largest singular values of C.

Remark 5: When the kernel functions needed to form \mathbf{K}_x and \mathbf{K}_y are not available, one may presume $\mathbf{K}_x := \sum_{m=1}^{M} \theta_m \mathbf{K}_m$ and $\mathbf{K}_y := \sum_{m=1}^{M} \delta_m \mathbf{K}_m$ for (25). Here, $\{\mathbf{K}_m\}_{m=1}^{M}$ are known kernel matrices for a preselected dictionary of kernels, while $\{\theta_m, \delta_m\}_{m=1}^{M}$ are unknown coefficients to be optimized along with the canonical vectors through (25). Such a data-driven approach is also known as multi-kernel

Algorithm 2: Graph Kernel Canonical Correlation Analysis.

- 1: **Input:** $\{\mathbf{x}_i\}_{i=1}^N, \{\mathbf{y}_i\}_{i=1}^N, \mathbf{W}, d, \gamma, \epsilon, \kappa_x(\cdot), \text{ and } \kappa_y(\cdot).$
- 2: Construct \mathbf{K}_x and \mathbf{K}_y using (24).
- 3: **Build** L_G using (7).
- 4: Perform SVD on C := UΣV[⊤] in (26c), where the diagonal elements of Σ are organized in descending order; U ∈ ℝ^{N×N}, V ∈ ℝ^{N×N}, and Σ ∈ ℝ^{N×N}.
- 5: **Extract** the first *d* columns of **U** and **V** to form $\bar{\mathbf{A}}^* \in \mathbb{R}^{N \times d}$ and $\bar{\mathbf{B}}^* \in \mathbb{R}^{N \times d}$, respectively.
- 6: Compute $\mathbf{A}^* = \mathbf{K}_x^{-1/2} (\mathbf{K}_x + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{A}}^*$ and $\mathbf{B}^* = \mathbf{K}_y^{-1/2} (\mathbf{K}_y + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{B}}^*$.
- 7: Output: A^* and B^* .

TABLE I COMPUTATIONAL COMPLEXITY COMPARISON

| gCCA | $\mathcal{O}(\min(D_x, D_y)N^2)$ |
|--------------|----------------------------------|
| CCA | $\mathcal{O}(D^2N)$ |
| gdCCA (dCCA) | $\mathcal{O}(DN^2)$ |
| gKCCA (KCCA) | $\mathcal{O}(\max(D, N)N^2)$ |

learning, and it has been broadly studied; see for example, [3], [23].

In terms of computational cost, we summarize the complexities of gCCA, gdCCA, gKCCA, CCA, dCCA, and KCCA in Table I, where $D := \max(D_x, D_y)$. Note that gCCA incurs higher computational cost than standard CCA, due to the extra multiplication term of $\mathbf{XL}_{\mathcal{G}}\mathbf{Y}^T$ in gCCA. If $N \ll D$, then gCCA in its present form is not feasible, or suboptimal even if the pseudo-inverse or Tikhonov regularization is employed, at computational complexity $\mathcal{O}(D^3)$. In this case, gdCCA is computationally more attractive since its complexity grows only linearly with D. In terms of gKCCA, when $D \gg N$, evaluating the kernel matrices dominates the computational complexity, giving rise to $\mathcal{O}(DN^2)$. When $D \ll N$, Steps 4 and 6 in Alg. 2 dominate the complexity, incurring complexity of $\mathcal{O}(N^3)$.

VI. NUMERICAL TESTS

To showcase the merits of our novel approaches, several classification experiments using real datasets are reported in this section. Classification accuracies of our proposed gCCA, gd-CCA and gKCCA are compared with competing alternatives.

A. Tests for gCCA

In this experiment, the AR face dataset [22], and the Extended Yale-B (EYB) face image dataset [21], were used to examine the classification performance of different schemes, including gCCA, CCA, graph (g) PCA [27], PCA, graph regularized multi-set (GrM) CCA [35], and the *k*-nearest neighbors (KNN) method.

The AR face database contains color face images of 100 individuals, each depicted in 26 images. These 26 images per person were taken under different lighting conditions,

occlusions and expressions. Each image was cropped and resized to 40×30 pixels, converted to grayscale image, and vectorized to obtain a $1,200 \times 1$ vector. The 1,200 features of each image were unevenly split into two views, where one view consists of the first 300 features collected in one column of $\mathbf{X}_0 \in \mathbb{R}^{300 \times 2,600}$ (2,600 columns for all the images) , while the remaining 900 features were used to form $\mathbf{Y}_0 \in \mathbb{R}^{900 \times 2,600}$. Suppose that $N_{\rm tr}$ columns were randomly drawn from 26 columns of \mathbf{X}_0 and \mathbf{Y}_0 that correspond to one person, to form the training data $\mathbf{X} \in \mathbb{R}^{300 \times 100N_{\text{tr}}}$ and $\mathbf{Y} \in \mathbb{R}^{900 \times 100N_{\text{tr}}}$, respectively. For the remaining $(26 - N_{\rm tr})$ columns of \mathbf{X}_0 associated with each person, half of them were used for tuning the hyperparameters, and the other half for testing, which were collected in $\mathbf{X}_{tu} \in \mathbb{R}^{300 \times 100(13 - 0.5N_{tr})}$ and $\mathbf{X}_{te} \in \mathbb{R}^{300 \times 100(13 - 0.5N_{tr})}$ accordingly. Here, we consider the scenario where only one view, namely \mathbf{X}_{te} , is available in the testing phase, which is of practical importance when one only has partial information about the testing images.

The EYB database consists of frontal face images of 38 individuals, each of which has around 65 color images of 192 × 168 pixels. All images were resized to 30 × 20 pixels and converted to grayscale before being vectorized to obtain a 600 × 1 vector. Then, the vector associated with each image was split into two subvectors (views) with $D_x = 250$ and $D_y = 350$. For each individual, $N_{\rm tr}$ images were randomly selected and the corresponding two views were used to construct the training datasets $\mathbf{X} \in \mathbb{R}^{D_x \times 38N_{\rm tr}}$ and $\mathbf{Y} \in \mathbb{R}^{D_y \times 38N_{\rm tr}}$. Among the remaining images, $(30 - 0.5N_{\rm tr})$ images per individual were used for tuning dataset $\mathbf{X}_{\rm tu} \in \mathbb{R}^{D_x \times 38(30 - 0.5N_{\rm tr})}$ and another $(30 - 0.5N_{\rm tr})$ for testing dataset $\mathbf{X}_{\rm te} \in \mathbb{R}^{D_x \times 38(30 - 0.5N_{\rm tr})}$, after following a similar process to build \mathbf{X} .

Letting $N := 100N_{tr}$ for the AR data experiment ($N := 38N_{tr}$ for EYB), we collected all common sources $\{\mathbf{s}_i\}_{i=1}^N$ into matrix **S**, which was constructed using the training data as follows: $\mathbf{S} := [\mathbf{X}^\top \mathbf{Y}^\top]^\top = [\mathbf{s}_1 \cdots \mathbf{s}_N]$. Based on **S**, matrix **W** was formed to have (i, j)-th entry given by

$$w_{ij} := \begin{cases} \frac{\mathbf{s}_i^{\top} \mathbf{s}_j}{\|\mathbf{s}_i\|_2 \|\mathbf{s}_j\|_2} & \mathbf{s}_i \in \mathcal{N}_k(\mathbf{s}_j) \text{ or } \mathbf{s}_j \in \mathcal{N}_k(\mathbf{s}_i) \\ 0 & \text{ otherwise} \end{cases}$$
(27)

for i, j = 1, 2, ..., N, where $\mathcal{N}_k(\mathbf{s}_j)$ denotes the set of the *k*-nearest neighbors of \mathbf{s}_j that belong to the same class (person) in **S**. In this experiment, $k = N_{tr} - 1$ was kept fixed.

In this experiment, 30 Monte Carlo (MC) simulations were run to assess the classification performance of gCCA, standard CCA, GrMCCA, gPCA, PCA, and KNN on the AR face dataset, as well as the EYB dataset. For fairness, the weight matrix W in (27) was used for gPCA. The classification accuracy is defined as the ratio between the number of correctly classified images and the total number of images tested. For gCCA, CCA, GrMCCA, gPCA, and PCA, 50 (100) canonical vectors for the AR (EYB) face dataset were found to obtain the low-dimensional representations of testing data, which were subsequently classified through the 10-nearest neighbors algorithm based on the Euclidean distance metric. The hyper-parameters in gCCA, gPCA, and GrMCCA were tuned among 30 logarithmically-spaced



Fig. 1. Classification accuracy of gCCA on the AR face dataset [22].



Fig. 2. Classification accuracy of gCCA on the EYB dataset [21].

values between 10^{-3} and 10^{3} to maximize the classification accuracies on 'tuning set' of images.

Figures 1 and 2 depict the classification accuracies of gCCA, CCA, GrMCCA, gPCA, PCA, and KNN on the AR data, and the EYB data, respectively, for a varying number of training samples. It is evident that the accuracies of all simulated schemes improve as $N_{\rm tr}$ grows, and our proposed gCCA outperforms alternatives for $N_{\rm tr} \ge 10$. This corroborates that incorporating the source graph that encodes dependencies among common sources, pays off.

B. Tests for gdCCA

The second experiment evaluates the capability of gdCCA for classification using again the AR face dataset and the EYB dataset. Per MC run on the AR face dataset, we collected all images of 10 randomly sampled people. For each selected person, $N_{\rm tr}$, $(13 - 0.5N_{\rm tr})$, and $(13 - 0.5N_{\rm tr})$ images were randomly drawn for training, tunning, and testing, respectively. In the training phase, each image was first converted to a grayscale image, resized to 80×60 pixels, and subsequently lexicographically ordered to obtain a $4,800 \times 1$



Fig. 3. Classification accuracy of gdCCA using dataset [22].

vector. To create the two views, this vector was partitioned into two subvectors of size $D_x = 1,000$ for $\mathbf{X} \in \mathbb{R}^{D_x \times 10N_{\text{tr}}}$ and of size $D_y = 3,800$ for $\mathbf{Y} \in \mathbb{R}^{D_y \times 10N_{\text{tr}}}$. Similarly, the training data $\mathbf{X}_{\text{tu}} \in \mathbb{R}^{D_x \times 10(13 - 0.5N_{\text{tr}})}$ and testing data $\mathbf{X}_{\text{te}} \in \mathbb{R}^{D_x \times 10(13 - 0.5N_{\text{tr}})}$ were generated.

Per realization on the EYB dataset, images of 10 individuals were randomly selected, and the two-view data $\mathbf{X} \in \mathbb{R}^{D_x \times 10N_{\text{tr}}}$ and $\mathbf{Y} \in \mathbb{R}^{D_y \times 10N_{\text{tr}}}$ were generated using the same procedure described for the AR data, except for $D_x = 1,000$ and $D_y =$ 7,000. For both the tuning data $\mathbf{X}_{\text{tu}} \in \mathbb{R}^{D_x \times 10(30-0.5N_{\text{tr}})}$ and the testing data $\mathbf{X}_{\text{te}} \in \mathbb{R}^{D_x \times 10(30-0.5N_{\text{tr}})}$, a number of $(30 - 0.5N_{\text{tr}})$ images were randomly chosen per person.

The two-view data in the training phase formed $\mathbf{S} = [\mathbf{X}^{\top} \mathbf{Y}^{\top}]^{\top}$ and were further used to build \mathbf{W} as in (27). For fairness, graph dual (gd) PCA [27] was tested with the same \mathbf{W} as in gdCCA. Moreover, the two associated graph adjacency matrices in Laplacian regularized (Lr) CCA [7] were constructed via (27) after substituting \mathbf{S} by \mathbf{X} and \mathbf{Y} , respectively. We tuned the hyper-parameters in gdCCA, dual (d) CCA, LrCCA and gdPCA among 30 logarithmically spaced values between 10^{-3} and 10^3 to maximize the classification accuracy on data \mathbf{X}_{tu} . Here, dCCA was implemented by gdCCA after assigning $\gamma = 0$. In gdCCA, dCCA, LrCCA, gdPCA and dPCA [27], 20 and 100 projection vectors were used for obtaining lower-dimensional representations of \mathbf{X}_{te} for AR data and EYB data, respectively. Then, the KNN rule with K = 10 was applied to carry out the classification tasks.

Figures 3 and 4 present the averaged classification accuracies of gdCCA, dCCA, LrCCA, gdPCA, dPCA, and KNN for a varying number of training images per person over 30 MC realizations. Clearly, our gdCCA enjoys the best classification performance among all simulated schemes for different training samples.

There are two hyper-parameters, namely γ and ϵ in gdCCA. To understand how the hyper-parameters influence the classification performance, the gdCCA was simulated on the AR face dataset for a range of γ and ϵ values. For each person, 17 (9) images were employed for training (testing). Fig. 5 plots the averaged classification accuracies over 30 MC runs, with γ varying from 10^{-3} to 10^3 and ϵ from 10^{-5} to 10^3 . For small γ



Fig. 4. Classification accuracy of gdCCA using dataset [21].



Fig. 5. Classification accuracy of gdCCA versus γ and ϵ .

values, the performance of gdCCA with small ϵ values outperforms that using large ϵ values. This is because with small γ , gdCCA approximates dCCA, and the Tikhonov regularization with excessively large ϵ values dominates the term for promoting uncorrelatedness between canonical variables. When ϵ is small, with γ increasing, the classification accuracy gradually increases by progressively exploiting the graph information, but subsequently decreases due to discarding the maximization of canonical correlations. Those observations confirm the assertion that with properly selected and nonzero γ and ϵ values, the performance of gdCCA reaches the best, in which case both maximizing the canonical correlations and exploiting the graph knowledge are in effect.

C. Tests for gKCCA

This last experiment assesses gKCCA for classification using the MNIST dataset.² There are 10 classes of handwritten 28×28 grayscale digit images in the MNIST, and each class (digit) consists of 7,000 images. Per MC run, 5 classes of images were randomly sampled for classification. For each selected class, $N_{\rm tr}$, $0.5N_{\rm tr}$, and $0.5N_{\rm tr}$ images

²Downloaded from http://yann.lecun.com/exdb/mnist/.

were randomly sampled for training, parameter tuning, and testing, respectively. The two-view data were created as follows. The images were first resized to 20×20 pixels, followed by vectorization. Each vector was split to 2 subvectors of sizes D_x and $D_y = 400 - D_x$ for the two views. The first/second view of training data was denoted by training dataset $\mathbf{X} \in \mathbb{R}^{D_x \times 5N_{\text{tr}}} / \mathbf{Y} \in \mathbb{R}^{D_y \times 5N_{\text{tr}}}$. The tuning/testing dataset $\mathbf{X}_{\text{tu}} / \mathbf{X}_{\text{te}}$ were the first views of tuning/testing images.

Gaussian kernels were used for X, Y, and the common source $\mathbf{S} := [\mathbf{X}^{\top} \mathbf{Y}^{\top}]^{\top}$, whose bandwidth parameters were set as the medians of the corresponding Euclidean distances. The idea to generate the W in Section VI-A was adopted and adjusted for constructing the graph adjacency matrix, which was also denoted by W for notational simplicity. Obviously, the similarity between two sources in S can not be measured by the linear correlation coefficient, which instead can be represented by a corresponding element in the kernel matrix of S, namely \mathbf{K}_s . Specifically,

$$w_{ij} := \begin{cases} \mathbf{K}_s(i,j) & \mathbf{s}_i \in \mathcal{M}_{k_1}(\mathbf{s}_j) \text{ or } \mathbf{s}_j \in \mathcal{M}_{k_1}(\mathbf{s}_i) \\ 0 & \text{otherwise} \end{cases}$$
(28)

for $i, j = 1, 2, ..., 5N_{tr}$, where s_i denotes the *i*-th source (*i*-th column) in S, and $\mathcal{M}_{k_1}(\mathbf{s}_i)$ is the set containing the k_1 -nearest neighbors of s_i from the same class. In the simulations of this subsection, $k_1 = N_{tr} - 1$. Further, graph kernel (gK) PCA [27] was simulated with the same W as in gKCCA. The graph Laplacian regularized (Lr) KCCA [7] was associated with two graph adjacency matrices, which were obtained by (28) after substituting \mathbf{K}_s with \mathbf{K}_x and \mathbf{K}_y accordingly. For fairness, all the kernel-based methods, namely gKCCA, KCCA, LrKCCA, gKPCA, and KPCA, shared the same kernel \mathbf{K}_x (and \mathbf{K}_y). When implementing the CCA-based and PCA-based subspace methods, 20 projection vectors were used for classification using the K-NN algorithm with K = 10. The hyper-parameters of gKCCA, KCCA, gdCCA, dCCA, LrKCCA, LrCCA, gKPCA, and gdPCA, were selected from 30 logarithmically spaced values between 10^{-3} and 10^{3} . For each algorithm, the parameters were selected with the best classification accuracy on the tuning dataset X_{tu} . In the following tests, the classification performance of all aforementioned algorithms was achieved after running 30 independent realizations.

In Fig. 6, the classification accuracies of simulated schemes for a variable number of training samples are reported, with $D_x = 120$ and $D_y = 280$. The plots validate the advantage of our gKCCA relative to the other 10 methods. Moreover, with extra training samples becoming available, the performance of all simulated schemes improves. Figure 7 depicts the classification accuracies of all methods for different D_x values, with $N_{tr} = 30$ kept fixed. It is clear that gKCCA outperforms alternatives under different vector splittings. Meanwhile, with D_x decreasing, it becomes more challenging to classify the testing data, so the classification accuracies of all schemes decrease. Interestingly, the performance gap between gKCCA and the others widens for smaller D_x values.



Fig. 6. Classification accuracy of gKCCA versus $N_{\rm tr}$.



Fig. 7. Classification accuracy of gKCCA versus D_x .

VII. CONCLUSION

Graph regularized CCA, dual CCA, as well as kernel CCA methods were revisited in this paper to exploit hidden lowdimensional common structures from two-view data of the same sources. Distinguishing itself from prior CCA contributions, our gCCA framework leverages additional information to improve the low-dimensional approximations through the canonical variables, by embedding the hidden common sources in a graph and invoking this graph prior knowledge as a CCA regularizer. As such, canonical pairs that are able to capture the structural information between data vectors can be revealed. In certain practical setups where the number of data samples is small relative to the data vector dimensionality, our gCCA is not directly applicable, or leads to suboptimal performance and incurs high computational complexity. To bypass this, the dual model of gCCA, namely gdCCA, is put forth. To further account for nonlinear data dependencies, the graph kernel CCA is developed. Numerical tests on several real-world datasets are presented to demonstrate the merits of the novel approaches.

This paper opens up several intriguing directions for future research. To start, evaluating analytically the performance of proposed schemes using (possibly) concentration inequality bounds is pertinent. Developing data-driven approaches to select the appropriate kernels (graphs) from a given or constructed dictionary of kernels (graphs) [28] is also meaningful. To endow the proposed gCCA algorithms with scalability, distributed and online implementations are well-motivated for handling largescale and/or high-dimensional streaming data. Generalizing our gCCA models to unpaired or multi-view datasets constitutes another interesting direction.

APPENDIX A

A. Proof of Theorem 1

Letting

$$ar{\mathbf{U}} := \mathbf{\Sigma}_x^{1/2} \mathbf{U} \in \mathbb{R}^{D_x imes d}, \quad ext{and} \quad ar{\mathbf{V}} := \mathbf{\Sigma}_y^{1/2} \mathbf{V} \in \mathbb{R}^{D_y imes d}$$

the objective function (12a) can be rewritten as

$$\operatorname{Tr}(\bar{\mathbf{U}}^{\top}\mathbf{C}\bar{\mathbf{V}}) := \operatorname{Tr}(\bar{\mathbf{U}}^{\top}\boldsymbol{\Sigma}_{x}^{-1/2}(\boldsymbol{\Sigma}_{xy} - \gamma\mathbf{X}\mathbf{L}_{\mathcal{G}}\mathbf{Y}^{\top})\boldsymbol{\Sigma}_{y}^{-1/2}\bar{\mathbf{V}})$$

and problem (12) boils down to

$$\max_{\bar{\mathbf{U}}, \bar{\mathbf{V}}} \operatorname{Tr}(\bar{\mathbf{U}}^{\top} \mathbf{C} \bar{\mathbf{V}})$$
(29a)

s. to $\bar{\mathbf{U}}^{\top}\bar{\mathbf{U}} = \mathbf{I}$, and $\bar{\mathbf{V}}^{\top}\bar{\mathbf{V}} = \mathbf{I}$. (29b)

Clearly, problem (29) is a typical truncated SVD formulation. So the columns of the optimal solutions $\overline{\mathbf{U}}^*$ and $\overline{\mathbf{V}}^*$ collect the *d* left and right singular vectors of **C** associated with the first *d* largest singular values, respectively.

Once having computed $\overline{\mathbf{U}}^*$ and $\overline{\mathbf{V}}^*$, the optimal solutions \mathbf{U}^* and \mathbf{V}^* to problem (12) are obtained as $\mathbf{U}^* := \Sigma_x^{-1/2} \overline{\mathbf{U}}^*$ and $\mathbf{V}^* := \Sigma_y^{-1/2} \overline{\mathbf{V}}^*$. Moreover, the maximal value of (12a) becomes the sum of the first *d* largest singular values of \mathbf{C} .

APPENDIX B

B. Proof of Theorem 2

Upon defining

$$ar{\mathbf{A}} := (\mathbf{K}_{\mathbf{x}} + \epsilon \mathbf{I})^{1/2} \mathbf{K}_{x}^{1/2} \mathbf{A}$$

 $ar{\mathbf{B}} := (\mathbf{K}_{\mathbf{v}} + \epsilon \mathbf{I})^{1/2} \mathbf{K}_{y}^{1/2} \mathbf{B}$

problem (25) can be rewritten as

$$\begin{split} (\bar{\mathbf{A}}^*, \, \bar{\mathbf{B}}^*) &:= \arg \max_{\bar{\mathbf{A}}, \, \bar{\mathbf{B}}} \ \mathrm{Tr}(\bar{\mathbf{A}}^\top \mathbf{C} \bar{\mathbf{B}}) \\ &\text{s. to } \ \bar{\mathbf{A}}^\top \bar{\mathbf{A}} = \mathbf{I}, \text{ and } \bar{\mathbf{B}}^\top \bar{\mathbf{B}} = \mathbf{I}. \end{split}$$

Using the results in Appendix A, one readily concludes that the columns of optimizers $\bar{\mathbf{A}}^*$, $\bar{\mathbf{B}}^*$ consist of the *d* left and right singular vectors of **C** associated with the first *d* largest singular values, respectively, which leads to

$$\mathbf{A}^* = \mathbf{K}_x^{-1/2} (\mathbf{K}_x + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{A}}^*$$
$$\mathbf{B}^* = \mathbf{K}_y^{-1/2} (\mathbf{K}_y + \epsilon \mathbf{I})^{-1/2} \bar{\mathbf{B}}^*.$$

Likewise, the maximal value of (25a) is given by the sum of the d largest singular values of **C**.

REFERENCES

- G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *Proc. Int. Conf. Mach. Learn.*, Atlanta, GA, USA, Jun. 2013, pp. 1247–1255.
- [2] N. Aronszajn, "Theory of reproducing kernels," *Trans. Amer. Math. Soc.*, vol. 68, no. 3, pp. 337–404, 1950.
- [3] F. R. Bach, G. R. Lanckriet, and M. I. Jordan, "Multiple kernel learning, conic duality, and the SMO algorithm," in *Proc. Int. Conf. Mach. Learn.*, New York, NY, USA, Jul. 2004, pp. 1–8.
- [4] M. Belkin and P. Niyogi, "Laplacian Eigenmaps for dimensionality reduction and data representation," *Neural Comput.*, vol. 15, no. 6, pp. 1373– 1396, Jun. 2003.
- [5] M. Belkin, P. Niyogi, and V. Sindhwani, "Manifold regularization: A geometric framework for learning from labeled and unlabeled examples," *J. Mach. Learn. Res.*, vol. 7, pp. 2399–2434, Nov. 2006.
- [6] D. P. Bertsekas, *Nonlinear Programming*, 2nd ed. Belmont, MA, USA: Athena Scientific, 1999.
- [7] M. B. Blaschko, J. A. Shelton, A. Bartels, C. H. Lampert, and A. Gretton, "Semi-supervised kernel canonical correlation analysis with application to human fMRI," *Pattern Recognit. Lett.*, vol. 32, no. 11, pp. 1572–1583, Aug. 2011.
- [8] J. A. Brown, J. D. Rudie, A. Bandrowski, J. D. van Horn, and S. Y. Bookheimer, "The UCLA multimodal connectivity database: A web-based platform for brain connectivity matrix sharing and analysis," *Front Neuroinform.*, vol. 6, Nov. 2012, Art. no. 28.
- [9] J. Chen and I. D. Schizas, "Online distributed sparsity-aware canonical correlation analysis," *IEEE Trans. Signal Process.*, vol. 64, no. 3, pp. 688– 703, Feb. 2016.
- [10] J. Chen, G. Wang, and G. B. Giannakis, "Nonlinear dimensionality reduction for discriminative analytics of multiple datasets," *IEEE Trans. Signal Process.*, May 2018, submitted. [Online]. Available: https://arxiv.org/abs/1805.05502
- [11] J. Chen, G. Wang, Y. Shen, and G. B. Giannakis, "Canonical correlation analysis with common graph priors," in *Proc. Statist. Signal Process. Workshop*, Freiburg, Germany, Jun. 2018, pp. 488–492.
- [12] J. Chen, A. Malhotra, and I. D. Schizas, "Data-driven sensors clustering and filtering for communication efficient field reconstruction," *Signal Process.*, vol. 133, pp. 156–168, Apr. 2017.
- [13] N. M. Correa, T. Adali, Y.-O. Li, and V. D. Calhoun, "Canonical correlation analysis for data fusion and group inferences," *IEEE Signal Process. Mag.*, vol. 27, no. 4, pp. 39–50, Jun. 2010.
- [14] G. B. Giannakis, Y. Shen, and G. V. Karanikolas, "Topology identification and learning over graphs: Accounting for nonlinearities and dynamics," *Proc. IEEE*, vol. 106, no. 5, pp. 787–807, May 2018.
- [15] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical correlation analysis: An overview with application to learning methods," *Neural Comput.*, vol. 16, no. 12, pp. 2639–2664, Dec. 2004.
- [16] H. Hotelling, "Relations between two sets of variates," *Biometrika*, vol. 28, no. 3/4, pp. 321–377, Dec. 1936.
- [17] B. Jiang, C. Ding, and J. Tang, "Graph-Laplacian PCA: Closed-form solution and robustness," in *Proc. Int. Conf. Comput. Vis. Pattern Recognit.*, Portland, OR, USA, Jun. 2013, pp. 3492–3498.
- [18] I. T. Jolliffe, Principal Component Analysis. Hoboken, NJ, USA: Wiley, 2002.
- [19] F. R. S. Karl Pearson, "LIII. On lines and planes of closest fit to systems of points in space," *London Edinburgh Dublin Philos. Mag. J. Sci.*, vol. 2, no. 11, pp. 559–572, 1901.
- [20] V. Kekatos, G. Wang, H. Zhu, and G. B. Giannakis, "PSSE redux: Convex relaxation, decentralized, robust, and dynamic approaches," in *Advances* in *Electric Power and Energy; Power Systems Engineering*, M. El-Hawary, Ed. Hoboken, NJ, USA: Wiley-Blackwell, 2018.
- [21] K. C. Lee, J. Ho, and D. J. Kriegman, "Acquiring linear subspaces for face recognition under variable lighting," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 5, pp. 684–698, May 2005.
- [22] A. R. Martinez and R. Benavente, "The AR face database, 1998," Computer Vision Center, Barcelona, Spain, Tech. Rep. #24, 2007, vol. 3, p. 5.
- [23] D. Romero, M. Ma, and G. B. Giannakis, "Kernel-based reconstruction of graph signals," *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 764–778, Feb. 2017.
- [24] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, vol. 290, no. 5500, pp. 2323–2326, Dec. 2000.
- [25] N. Shahid, N. Perraudin, V. Kalofolias, G. Puy, and P. Vandergheynst, "Fast robust PCA on graphs," *IEEE J. Sel. Top. Signal Process.*, vol. 10, no. 4, pp. 740–756, Feb. 2016.

- [26] F. Shang, L. Jiao, and F. Wang, "Graph dual regularization non-negative matrix factorization for co-clustering," *Pattern Recognit.*, vol. 45, no. 6, pp. 2237–2250, Jun. 2012.
- [27] Y. Shen, P. Traganitis, and G. B. Giannakis, "Nonlinear dimensionality reduction on graphs," in *Proc. IEEE Int. Workshop Comput. Adv. Multi*sensor Adaptive Process., Curacao, Dutch Antilles, Dec. 2017, pp. 1–5.
- [28] Y. Shen, P. Traganitis, and G. B. Giannakis, "Graph-adaptive nonlinear dimensionality reduction," *IEEE Trans. Signal Process.*, Mar. 2018, submitted to be published. [Online]. Available: https://arxiv.org/pdf/1801.09390.pdf
- [29] Y. Shen, B. Baingana, and G. B. Giannakis, "Kernel-based structural equation models for topology identification of directed networks," *IEEE Trans. Signal Process.*, vol. 65, no. 10, pp. 2503–2516, May 2017.
- [30] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending highdimensional data analysis to networks and other irregular domains," *IEEE Signal Process. Mag.*, vol. 30, no. 3, pp. 83–98, May 2013.
- [31] A. J. Smola and R. Kondor, "Kernels and regularization on graphs," in *Learning Theory and Kernel Machines*. Berlin, Germany: Springer, 2003, pp. 144–158.
- [32] S. VanVaerenbergh, J. Vía, and I. Santamaría, "Blind identification of SIMO wiener systems based on kernel canonical correlation analysis," *IEEE Trans. Signal Process.*, vol. 61, no. 9, pp. 2219–2230, May 2013.
- [33] D. M. Witten, R. Tibshirani, and T. Hastie, "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis," *Biostatistics*, vol. 10, no. 3, pp. 515–534, Apr. 2009.
- [34] Y. Yamanishi, J.-P. Vert, A. Nakaya, and M. Kanehisa, "Extraction of correlated gene clusters from multiple genomic data by generalized kernel canonical correlation analysis," *Bioinformatics*, vol. 19, no. 1, pp. i323– i330, Jul. 2003.
- [35] Y. Yuan and Q. Sun, "Graph regularized multiset canonical correlations with applications to joint feature extraction," *Pattern Recognit.*, vol. 47, no. 12, pp. 3907–3919, Dec. 2014.



Jia Chen received the B.S. degree from the Southwest Jiaotong University, Chengdu, China, in 2009, and the M.S. degree from the University of Electronic Science and Technology of China, Chengdu, China, in 2012, and the Ph.D. degree from the University of Texas at Arlington, Arlington, TX, USA, in 2016, all in electrical engineering. She is currently a Postdoctoral Associate with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA. Her research interests include signal processing, data analytics, and machine learning.



Gang Wang (M'18) received the B.Eng. degree in electrical engineering and automation from the Beijing Institute of Technology, Beijing, China, in 2011, and the Ph.D. degree in electrical and computer engineering from the University of Minnesota, Minneapolis, MN, USA, in 2018. He is currently a Postdoctoral Associate with the

Department of Electrical and Computer Engineering at University of Minnesota. His research interests focus on the areas of statistical signal processing, stochastic and nonconvex optimization with applica-

tions to autonomous energy grids, and deep learning. He was a recipient of the National Scholarship (2014), a Guo Rui Scholarship (2017), and an Innovation Scholarship (first place in 2017), all from China, as well as a Best Student Paper Award at the 2017 European Signal Processing Conference.



Yanning Shen (S'13) received the B.Sc. and M.Sc. degrees in electrical engineering from the University of Electronic Science and Technology of China, Chengdu, China, in 2011 and 2014, respectively. Since September 2014, she has been working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN, USA. Her research interests include signal processing on graphs, network science, and machine learning. She was a Best Student Paper Award finalist of the 2017 IEEE International Work-

shop on Computational Advances in Multisensor Adaptive Processing. She was selected to participate in the 2017 Rising Stars in EECS Workshop at Stanford University, and she was a recipient of the UMN Doctoral Dissertation Fellowship in 2018.



Georgios B. Giannakis (F'97) received the Diploma degree in electrical engineering from the National Technical University of Athens, Athens, Greece, in 1981. From 1982 to 1986, he was with the University of the Southern California, Los Angeles, CA, USA, where he received the M.Sc. degree in electrical engineering in 1983, the M.Sc. degree in mathematics, and the Ph.D. degree in electrical engineering both in 1986.

He was with the University of Virginia from 1987 to 1998 and since 1999, he has been a Professor with

the University of Minnesota, Minneapolis, MN, USA, where he holds an Endowed Chair in wireless telecommunications, the University of Minnesota McKnight Presidential Chair in ECE, and the Director of the Digital Technology Center. His general interests include the areas of communications, networking, and statistical signal processing-subjects on which he has authored or co-authored more than 400 journal papers, 700 conference papers, 25 book chapters, two edited books, and two research monographs (h-index 132). Current research focuses on learning from big data, wireless cognitive radios, and network science with applications to social, brain, and power networks with renewables. He is a fellow of the EURASIP, and was with the IEEE in a number of posts including that of a Distinguished Lecturer for the IEEE-SP Society. He is the (co-)inventor of 30 patents issued, and the (co-)recipient of nine best paper awards from the IEEE Signal Processing and Communications Societies including the G. Marconi Prize Paper Award in Wireless Communications. He is also a recipient of the Technical Achievement Awards from the SP Society (2000), the EURASIP (2005), a Young Faculty Teaching Award, the G. W. Taylor Award for Distinguished Research from the University of Minnesota, and the IEEE Fourier Technical Field Award (inaugural recipient in 2015).