

Scalable and Flexible Multiview Canonical Correlation Analysis

Xiao Fu, Kejun Huang, Mingyi Hong, Nicholas D. Sidiropoulos, and Anthony M.-C. So

Abstract— This paper considers generalized (multiview) canonical correlation analysis (GCCA) for large-scale datasets. A memory-efficient and computationally lightweight algorithmic framework is proposed for the classic MAX-VAR GCCA formulation as well as its variants. GCCA is gaining renewed interest in various applications such as speech processing and natural language processing. The classic MAX-VAR GCCA problem can be solved optimally via eigen-decomposition of a matrix that compounds the (whitened) correlation matrices of the views. However, this route can easily lead to memory explosion and a heavy computational burden when the size of the views becomes large. In addition, it was unclear how to promote per-specified structure (e.g. sparsity) on the canonical components sought, while structured components analysis is often desired in data analytics. In this work, we propose an alternating optimization (AO)-based algorithm to handle large-scale MAX-VAR GCCA as well as its variations that impose structure on the canonical components. The algorithm avoids instantiating the correlation matrices of the views and thus can achieve substantial saving in memory. It also maintains data sparsity, which can be exploited to alleviate the computational burden. Consequently, the proposed algorithm is highly scalable. Regularization such as sparsity promoting functions can be easily incorporated in the proposed computational framework. Convergence properties of the proposed algorithm are carefully studied. Simulations and large-scale word embedding tasks are employed to showcase the effectiveness of the proposed algorithm.

I. INTRODUCTION

Canonical Correlation Analysis (CCA) [1] produces low dimensional representations by finding common structure of two or more views corresponding to the same entities. A view contains high-dimensional representations of the entities in some domain – e.g., the text and audio representations corresponding to a given word can be considered as different views of this word. CCA is able to deal with views that have different dimensions, and this flexibility is very useful in data fusion, where one is interested in integrating information gathered from different domains. Multiview analysis finds numerous applications in signal processing and machine learning, such as blind source separation [2], [3], direction-of-arrival estimation [4], wireless channel equalization [5], regression [6], clustering [7], speech modeling and recognition [8], [9], and word embedding [10], to name a few. Classical

CCA was derived for the two-view case, but Generalized Canonical Correlation Analysis (GCCA) that aims at handling more than two views has a long history as well [11]. A typical application of GCCA is word embedding in natural language processing, where the vocabularies from different languages can be considered as multiple views of the same terms. Word embedding seeks low-dimensional representations of the terms that are well-aligned with human judgment, and applying GCCA to integrate multiple languages was shown to yield better embedding results relative to single-view analyses such as principle component analysis (PCA) [10].

Computationally, GCCA poses interesting and challenging optimization problems. Unlike the two-view case that admits an algebraically simple solution (via eigen-decomposition), GCCA is in general not easily solvable. Many prior works considered the GCCA problem with different cost functions [11], [12], while the proposed algorithms often can only extract a single canonical component and then find others through a deflation process, which is known to suffer from error propagation. Convergence properties of the GCCA algorithms are also largely under-investigated, since the GCCA formulations often involve non-convex constraints (e.g., manifold constraints) that complicate analysis. CCA and GCCA can also pose serious scalability challenges, since they involve auto- and cross-correlations of different views and a whitening stage [13]. These procedures can easily lead to memory explosion and require a large number of flops for computation. They also destroy the sparsity of the data, which is usually what one relies upon to deal with large-scale problems. In recent years, effort has been spent on solving these scalability issues, but the focus is mostly on the two-view case [13]–[15].

Among all different formulations of GCCA, there is a particular one that admits a conceptually simple solution, the so-called MAX-VAR GCCA [11], [12], [16]. MAX-VAR GCCA was first proposed in [11], and its solution amounts to finding the ‘directions’ aligned to those exhibiting maximum variance for a matrix aggregated from the (whitened) auto-correlations of the views. It can also be viewed as a problem of enforcing *identical* latent representations of different views as opposed to highly correlated ones, which is the more general goal of (G)CCA. The merit of MAX-VAR GCCA is that it can be solved via eigen-decomposition and finds all the canonical components simultaneously (i.e., no deflation is involved). In practice, MAX-VAR GCCA also demonstrates promising performance in various applications such as word embedding [10] and speech recognition [8]. On the other hand, MAX-VAR GCCA has the same scalability problem as the other GCCA formulations: It involves correlation matrices of

X. Fu, K. Huang and N.D. Sidiropoulos are with the Department of Electrical and Computer Engineering, University of Minnesota, Minneapolis, MN55455, e-mail (xfu,huang663,nikos)@umn.edu. M. Hong is with Department of Industrial and Manufacturing Systems Engineering, Iowa State University, Ames, Iowa 50011, (515) 294-4111, Email: mingyi@iastate.edu. Anthony M.-C. So is with the Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong, Shatin, N.T., Hong Kong, Email: manchoso@se.cuhk.edu.hk

different views and their inverses, which is prohibitive to even instantiate when the data dimension is large. The work in [10] provided a pragmatic way to circumvent this difficulty: PCA was first applied to each view to reduce the rank of the views, and then MAX-VAR GCCA was applied to the rank-truncated views. Such a procedure significantly reduces the number of parameters for characterizing the views and is feasible in terms of memory. However, truncating the rank of the views is prone to information loss, and thus leads to performance degradation.

Besides the basic (G)CCA formulations, *structured* CCA [17] that seeks canonical components with per-specified structure is often considered in applications. For example, sparse CCA is frequently seen in data analytics, for the purpose of discarding irrelevant features when performing CCA [18]–[20]. In multi-lingual word embedding [10], [14], [21], for example, it is known that there are many irrelevant features, such as stop words. In brain activation analysis, it is also believed that many voxels are irrelevant and should be somehow discarded [22], [23]. Gene analysis is another example [18]–[20]. Ideally, CCA seeks a few highly correlated latent components, and so it should naturally be able to identify and down-weight irrelevant features automatically. In practice, however, this ability is often impaired when correlations cannot be reliably estimated, when one only has access to relatively few and/or very noisy samples, or when there is model mismatch due to bad preprocessing (e.g., registration). In those cases, performing feature selection jointly with (G)CCA is well-motivated. However, introducing structure regularizations on the sought canonical components makes the optimization problem even harder – since many regularization terms like the sparsity-promoting regularizers are non-differentiable.

Contributions In this work, our interest lies in solving the MAX-VAR GCCA problem and its variants with structure-promoting regularizers when the dimensions of the views are large. Instead of truncating the rank of the views as in [10], we keep the data *intact* and devise a scalable algorithmic framework to handle the formulated problem. Specifically, our idea is to deal with problem using a two-block alternating optimization (AO) algorithm. Under the AO framework, the proposed algorithm alternates between a regularized least squares subproblem and an orthogonality-constrained subproblem. The merit of this framework is that correlation matrices of the views never need to be explicitly instantiated, and the inversion procedure is avoided. Consequently, the algorithm consumes significantly less memory compared to that required by the original solution using eigen-decomposition and is also very flexible in incorporating different structure-promoting regularizers.

On the theory side, convergence properties of the algorithm are also carefully studied: We first show that the proposed algorithm *globally* converges to a Karush-Kuhn-Tucker (KKT) point of the formulated problem under a variety of regularizers, even when the subproblems are solved in a grossly inexact manner. We also show that the optimality gap shrinks to at most $\mathcal{O}(1/r)$ after r iterations – i.e., at least a sublinear convergence rate can be guaranteed. In addition, we show that when there is no regularization, the proposed algorithm

solves the MAX-VAR problem to *global optimality* with a *linear* convergence rate. When the subproblems are not exactly solved, such a rate still holds with some loss in accuracy. The proposed algorithm is applied to synthetic data and a real large-scale word embedding problem, and promising results are observed.

Notation We use \mathbf{X} and \mathbf{x} to denote a matrix and a vector, respectively. $\mathbf{X}(m, :)$ and $\mathbf{X}(:, n)$ denote the m th row and the n th column of \mathbf{X} , respectively; in particular, $\mathbf{X}(:, \ell_1 : \ell_2)$ denotes a submatrix of \mathbf{X} consisting of the ℓ_1 – ℓ_2 th columns of \mathbf{X} (MATLAB notation). $\|\mathbf{X}\|_F$ and $\|\mathbf{X}\|_p$ for $p \geq 1$ denote the Frobenius norm and the matrix-induced p -norm, respectively. $\|\mathbf{X}\|_{p,1} = \sum_{i=1}^m \|\mathbf{X}(i, :)\|_p$ for $p \geq 1$ denotes the $\ell_p \ell_1$ -mixed norm of $\mathbf{X} \in \mathbb{R}^{m \times n}$. The superscripts “ T ”, “ \dagger ”, and “ -1 ” denote the matrix operators of transpose, pseudo-inverse and inverse, respectively. The operator $\langle \mathbf{X}, \mathbf{Y} \rangle$ denotes the inner product of \mathbf{X} and \mathbf{Y} (i.e., $\text{Tr}(\mathbf{X}^T \mathbf{Y})$).

II. BACKGROUND

The classic two-view CCA can be expressed as the following optimization problem [1]:

$$\begin{aligned} \min_{\mathbf{Q}_1, \mathbf{Q}_2} \quad & \|\mathbf{X}_1 \mathbf{Q}_1 - \mathbf{X}_2 \mathbf{Q}_2\|_F^2 \\ \text{s.t.} \quad & \mathbf{Q}_i^T (\mathbf{X}_i^T \mathbf{X}_i) \mathbf{Q}_i = \mathbf{I}, \quad i = 1, 2, \end{aligned} \quad (1)$$

where $\mathbf{X}_i \in \mathbb{R}^{L \times M_i}$, with its ℓ th row $\mathbf{X}_i(\ell, :)$ containing the i th view of the ℓ th data point, the columns of $\mathbf{Q}_i \in \mathbb{R}^{M_i \times K}$ correspond to the K canonical components of view \mathbf{X}_i , and $(1/L)\mathbf{X}_i^T \mathbf{X}_i$ serves as an estimate of the correlation of each view. Note that we are essentially maximizing the trace of the estimated cross-correlations between the views, i.e., $\text{Tr}(\mathbf{Q}_2^T \mathbf{X}_2^T \mathbf{X}_1 \mathbf{Q}_1)$. Thus, \mathbf{Q}_1 and \mathbf{Q}_2 can be considered as two linear operators that reduce the dimensionality of \mathbf{X}_1 and \mathbf{X}_2 , respectively, so that the reduced-dimension views, i.e., $\mathbf{X}_i \mathbf{Q}_i$ for $i = 1, 2$, are highly correlated. The constraints serve the purpose of normalization. Problem (1) can be solved exactly via an (generalized) eigen-decomposition, but this simple solution only applies to the two-view case. To analyze the case with more than two views, the so-called generalized CCA (GCCA) is often adopted [11]. Different cost functions of GCCA were proposed [11], [12], [24], [25], and arguably the most natural extension of CCA is the following [1]:

$$\begin{aligned} \min_{\{\mathbf{Q}_i\}_{i=1}^I} \quad & \sum_{i=1}^{I-1} \sum_{j=i+1}^I \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{X}_j \mathbf{Q}_j\|_F^2 \\ \text{s.t.} \quad & \mathbf{Q}_i^T (\mathbf{X}_i^T \mathbf{X}_i) \mathbf{Q}_i = \mathbf{I}, \quad i = 1, \dots, I, \end{aligned} \quad (2)$$

where I is the number of views. Unfortunately, Problem (2) does not admit an analytic solution. Another formulation of GCCA is more tractable: Instead of forcing pairwise similarity of the reduced-dimension views, one can seek a common latent representation of different views, i.e., [8], [10]–[12], [16]

$$\min_{\{\mathbf{Q}_i\}_{i=1}^I, \mathbf{G}} \sum_{i=1}^I (1/2) \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2, \quad (3)$$

where $\mathbf{G} \in \mathbb{R}^{L \times K}$ is the common latent representation of the different views. Conceptually, Problems (3) and (2) share

the same goal, i.e., find highly correlated reduced-dimension views. On the other hand, the solutions may be different since the $\mathbf{X}_i \mathbf{Q}_i$'s yielded by Problem (3) are not forced to be exactly orthogonal, but only approximately so. The upshot of Problem (3) is that it admits a *conceptually* simple algebraic solution, which, as we will show, has the potential to be scaled up to deal with very large problems. In this work, we will focus on Problem (3).

Problem (3) is referred to as the MAX-VAR formulation of GCCA since the optimal solution amounts to taking principal eigenvectors of a matrix aggregated from the correlation matrices of the views. To explain, let us first assume that \mathbf{X}_i has full column rank and marginalize \mathbf{Q}_i by letting $\mathbf{Q}_i = (\sqrt{L}) \mathbf{X}_i^\dagger \mathbf{G}$, where $\mathbf{X}_i^\dagger = (\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T$. By substituting it back to (3), we see that an optimal solution \mathbf{G}_{opt} can be obtained via solving the following:

$$\mathbf{G}_{\text{opt}} = \arg \max_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \text{Tr} \left(\mathbf{G}^T \left(\sum_{i=1}^I \mathbf{X}_i \mathbf{X}_i^\dagger \right) \mathbf{G} \right). \quad (4)$$

Let $\mathbf{M} = \sum_{i=1}^I \mathbf{X}_i \mathbf{X}_i^\dagger$. Then, an optimal solution is $\mathbf{G}_{\text{opt}} = \mathbf{U}_M(:, 1 : K)$, the first K principal eigenvectors of \mathbf{M} [26].

Although Problem (3) admits a seemingly easy solution, implementing it for large-scale data is prohibitive. The first difficulty lies in memory: As mentioned, instantiating $\mathbf{M} = \sum_{i=1}^I \mathbf{X}_i (\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T$ is not doable when L and M_i 's are large. The matrix \mathbf{M} is an $L \times L$ matrix. In applications like word embedding, L and M_i are the vocabulary size of a language and the number of features defining the terms, respectively, which can both easily exceed 100,000. This means that the memory for simply instantiating \mathbf{M} or $(\mathbf{X}_i^T \mathbf{X}_i)^{-1}$ can reach 75GB. In addition, even if the views \mathbf{X}_i are sparse, computing $(\mathbf{X}_i^T \mathbf{X}_i)^{-1}$ will create large dense matrices and make it difficult to exploit sparsity in the subsequent processing. To circumvent these difficulties, Rastogi *et al.* [10] proposed to first apply the singular value decomposition (SVD) to the views, i.e., $\text{svd}(\mathbf{X}_i) = \mathbf{U}_i \Sigma_i \mathbf{V}_i^T$, and then let

$$\hat{\mathbf{X}}_i = \mathbf{U}_i(:, 1 : P) \Sigma_i(1 : P, 1 : P) (\mathbf{V}_i(:, 1 : P))^T \approx \mathbf{X}_i,$$

where P is much smaller than M_i and L . This procedure enables one to represent the views with significantly fewer parameters, i.e., $(L + M_i + 1)P$ compared to LM_i , and allows the original eigen-decomposition based solution to MAX-VAR GCCA to be applied; see more details in [10]. The drawback, however, is also evident: The procedure essentially truncates the rank of the views significantly (since in practice the views almost always have full column-rank, i.e., $\text{rank}(\mathbf{X}_i) = M_i$), and rank-truncation may lose a lot of information. Therefore, it is much more appealing to deal with the original views directly. Another shortcoming of the above approach is that it is not flexible in incorporating regularizations on \mathbf{Q}_i , while structured \mathbf{Q}_i are often desired in different applications [18]–[20].

In this work, we provide an algorithmic framework that deals with the MAX-VAR GCCA problem and its variants with constraints on \mathbf{Q}_i . We aim at offering simple solutions that are memory-efficient, admit light per-iteration complexity,

and feature good convergence properties under certain mild conditions.

III. PROPOSED ALGORITHM

In this work, we consider a scalable and flexible algorithmic framework for handling MAX-VAR GCCA and some variants. Specifically, we consider the following formulation

$$\begin{aligned} \min_{\{\mathbf{Q}_i\}_{i=1}^I, \mathbf{G}} \quad & \sum_{i=1}^I (1/2) \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2 + \sum_{i=1}^I g_i(\mathbf{Q}_i), \\ \text{s.t.} \quad & \mathbf{G}^T \mathbf{G} = \mathbf{I}, \end{aligned} \quad (5)$$

where $g_i(\cdot)$ is a regularizer that imposes a certain structure on \mathbf{Q}_i . Popular regularizers are $g_i(\mathbf{Q}_i) = \mu_i \cdot \|\mathbf{Q}_i\|_F$, $g_i(\mathbf{Q}_i) = \mu_i \cdot \|\mathbf{Q}_i\|_{2,1}$ and $g_i(\mathbf{Q}_i) = \mu_i \cdot \|\mathbf{Q}_i\|_{1,1}$, where $\mu_i \geq 0$ is a regularization parameter for balancing the least squares fitting term and the regularization term. The first regularizer is commonly used for controlling the energy of the dimension-reducing matrix \mathbf{Q}_i , which also has an effect of improving the conditioning of the least squares problem. The latter two regularizers are used to select features automatically. To be specific, $g_i(\mathbf{Q}_i) = \|\mathbf{Q}_i\|_{2,1}$ promotes many rows of \mathbf{Q}_i to be zero (or approximately zero), and thus can suppress the impact of the corresponding columns (features) in \mathbf{X}_i – which is effectively feature selection. The function $g_i(\mathbf{Q}_i) = \|\mathbf{Q}_i\|_{1,1}$ also does feature selection, but different canonical components may use different features. In this section, we propose an algorithm that can deal with the regularized and the original version of MAX-VAR GCCA under a unified framework.

A. Proposed Algorithm: Alternating Optimization

To deal with Problem (5), our idea is alternating optimization; i.e., we solve two subproblems w.r.t. $\{\mathbf{Q}_i\}$ and \mathbf{G} , respectively. As will be seen, such a simple strategy will lead to highly scalable algorithms in terms of both memory and computational cost.

To begin with, let us consider the subproblem

$$\min_{\mathbf{Q}_i} (1/2) \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}^{(r)}\|_F^2 + g_i(\mathbf{Q}_i), \quad \forall i, \quad (6)$$

where $\mathbf{G}^{(r)}$ denotes the iterate of \mathbf{G} after the r th iteration. The above problem is a regularized least squares problem. When \mathbf{X}_i is large and sparse, many efficient algorithms can be considered to solve it. For example, the alternating direction method of multipliers (ADMM) [27] is frequently employed to handle Problem (6) in a scalable manner. However, ADMM is a primal-dual method that does not guarantee monotonic decrease of the objective value, which will prove useful in later convergence analysis. Hence, we propose to employ a simple proximal gradient method for handling Problem (6). By proximal gradient, we update \mathbf{Q}_i by the following update rule:

$$\mathbf{Q}_i^{(r+1)} \leftarrow \text{prox}_{g_i} \left(\mathbf{Q}_i^{(r)} - \alpha_i \nabla_{\mathbf{Q}_i} f \left(\mathbf{Q}_i^{(r)}, \mathbf{G}_i^{(r)} \right) \right), \quad (7)$$

where we define the proximal operator as

$$\text{prox}_g(\mathbf{y}) = \arg \min_{\mathbf{x}} 1/2 \|\mathbf{x} - \mathbf{y}\|_2^2 + g(\mathbf{x})$$

and use the notations $\mathbf{Q} = [\mathbf{Q}_1^T, \dots, \mathbf{Q}_I^T]^T$,

$$f(\mathbf{Q}, \mathbf{G}^{(r)}) = \sum_{i=1}^I \frac{1}{2} \left\| \mathbf{X}_i \mathbf{Q}_i - \mathbf{G}^{(r)} \right\|_F^2$$

and $\nabla_{\mathbf{Q}_i} f(\mathbf{Q}, \mathbf{G}_i^{(r)})$ as the partial gradient of the least squares part of the objective function w.r.t. \mathbf{Q}_i , i.e.,

$$\nabla_{\mathbf{Q}_i} f(\mathbf{Q}, \mathbf{G}_i^{(r)}) = \mathbf{X}_i^T \mathbf{X}_i \mathbf{Q}_i - \mathbf{X}_i^T \mathbf{G}^{(r)}. \quad (8)$$

For many functions $g_i(\cdot)$, the proximity operator in (7) has closed-form or lightweight solutions [28]. For example, for the regularization of interest such as $g_i(\mathbf{Q}_i) = \mu_i \|\mathbf{Q}_i\|_F$, the solution of Problem (7) is simply

$$\mathbf{Q}_i^{(r+1)} \leftarrow (\mathbf{X}_i^T \mathbf{X}_i + \mu_i \mathbf{I}) \mathbf{Q}_i^{(r)} - \mathbf{X}_i^T \mathbf{G}^{(r)}. \quad (9)$$

When $g_i(\mathbf{Q}_i) = \mu_i \|\mathbf{Q}_i\|_{2,1}$, the update rule becomes

$$\mathbf{Q}_i^{(r+1)}(m, :) \leftarrow \begin{cases} 0, & \|\mathbf{H}_i(m, :)\|_2 < \mu_i, \\ \left(1 - \frac{\mu_i}{\|\mathbf{H}_i(m, :)\|_2}\right) \|\mathbf{H}_i(m, :)\|_2, & \text{o.w.,} \end{cases} \quad (10)$$

where $\mathbf{H}_i = \mathbf{Q}_i^{(r)} - \alpha_i \nabla_{\mathbf{Q}_i} f(\mathbf{Q}^{(r)}, \mathbf{G}_i^{(r)})$. For $g_i(\mathbf{Q}_i) = \mu_i \|\mathbf{Q}_i\|_{1,1}$, the update rule is similar to that in (10), which is known as the *soft-thresholding operator*.

By updating \mathbf{Q}_i using the rule in (7) (for one time or several times), we obtain $\mathbf{Q}_i^{(r+1)}$. Next, we consider solving the subproblem w.r.t. \mathbf{G} when fixing $\{\mathbf{Q}_i\}_{i=1}^I$. Now we can drop the regularization term since it does not affect the cost value when \mathbf{Q}_i is fixed. Then, we have the following equivalence:

$$\begin{aligned} & \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \sum_{i=1}^I \frac{1}{2} \left\| \mathbf{X}_i \mathbf{Q}_i^{(r+1)} - \mathbf{G} \right\|_F^2 \\ & \Leftrightarrow \arg \max_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \text{Tr} \left(\mathbf{G}^T \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i^{(r+1)} / I \right). \end{aligned}$$

Therefore, an optimal solution of \mathbf{G} is as follows: Let $\mathbf{P} = \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i^{(r+1)}$. Then, we have

$$\mathbf{G}^{(r+1)} \leftarrow \mathbf{U}_P \mathbf{V}_P^T,$$

where $\mathbf{U}_P \Sigma_P \mathbf{V}_P^T = \text{svd}(\mathbf{P}, \text{'econ'})$, and $\text{svd}(\cdot, \text{'econ'})$ denotes the economy-size SVD that produces $\mathbf{U}_P \in \mathbb{R}^{L \times K}$, $\Sigma_P \in \mathbb{R}^{K \times K}$ and $\mathbf{V}_P^T \in \mathbb{R}^{K \times K}$. The above update is optimal in terms of solving the subproblem. In practice, one may also combine the knowledge of the previous iterate $\mathbf{G}^{(r)}$ and let

$$\mathbf{P} = \gamma \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i^{(r+1)} / I + (1 - \gamma) \mathbf{G}^{(r)}, \quad (11)$$

where $\gamma \in (0, 1]$. Such a slight modification of forming \mathbf{P} does not increase the complexity and is very easy to implement. However, such a simple combination helps establish nice convergence properties of the algorithm, as we will see in the next section.

The algorithm is summarized in Algorithm 1, which we call the alternating optimization-based MAX-VAR GCCA (AltMaxVar). As one can see, the algorithm does not instantiate any large dense matrix during the procedure and thus is highly efficient in terms of memory. Also, the procedure does

not destroy sparsity of the data, and thus the computational burden is light when the data is sparse – which is often the case in large-scale learning applications. Detailed complexity analysis will be presented in the next subsection.

Algorithm 1: AltMaxVar

```

input :  $\{\mathbf{X}_i, \mu_i\}_{i=1}^I; K; T; (\{\mathbf{Q}_i^{(0)}\}_{i=1}^I, \mathbf{G}^{(0)})$ .
1  $r \leftarrow 0$ ;
2 repeat
3    $t \leftarrow 0$ ;
4    $\mathbf{E}_i^{(t)} \leftarrow \mathbf{Q}_i^{(r)}$  for  $i = 1, \dots, I$ ;
5   while  $t < T$  and convergence not reached do
6     for all  $i$ , update
7        $\mathbf{E}_i^{(t+1)} \leftarrow \text{prox}_{g_i}(\mathbf{E}_i^{(t)} - \alpha_i \nabla f_{\mathbf{Q}_i}(\mathbf{E}_i^{(t)}; \mathbf{G}_i^{(r)}))$ 
8       where  $\nabla f_{\mathbf{Q}_i}(\mathbf{E}_i^{(t)}; \mathbf{G}_i^{(r)}) = -\mathbf{X}_i^T \mathbf{G}^{(r)} + \mathbf{X}_i^T \mathbf{X}_i \mathbf{E}_i^{(t)}$ ;
9        $t \leftarrow t + 1$ ;
10    end
11     $\mathbf{Q}_i^{(r+1)} \leftarrow \mathbf{E}_i^{(t)}$ ;
12     $\mathbf{P} \leftarrow \gamma \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i^{(r+1)} / I + (1 - \gamma) \mathbf{G}^{(r)}$ ;
13     $\mathbf{U}_P \mathbf{D}_P \mathbf{V}_P^T \leftarrow \text{svd}(\mathbf{P}, \text{'econ'})$ ;
14     $\mathbf{G}^{(r+1)} \leftarrow \mathbf{U}_P \mathbf{V}_P^T$ ;
15     $r \leftarrow r + 1$ ;
16 until Some stopping criterion is reached;
output:  $\{\mathbf{Q}_i^{(r)}\}_{i=1}^I, \mathbf{G}^{(r)}$ 

```

B. Computational and Memory Complexities

The update rule in (7) inherits the good features from the proximal gradient (PG) method. First, there is no “heavy computation” if the views \mathbf{X}_i for $i = 1, \dots, I$ are sparse. Specifically, the major computation in the update rule of (7) is computing the partial gradient of the smooth part of the cost function, i.e., $\nabla_{\mathbf{Q}_i} f(\mathbf{Q}_i, \mathbf{G}_i^{(r)})$. To this end, $\mathbf{X}_i \mathbf{Q}_i$ should be calculated first, since if \mathbf{X}_i is sparse, this matrix multiplication step has a complexity order of $\mathcal{O}(\text{nnz}(\mathbf{X}_i) \cdot K)$ flops, where $\text{nnz}(\cdot)$ counts the number of non-zeros. The next multiplication, i.e., $\mathbf{X}_i^T (\mathbf{X}_i \mathbf{Q}_i)$, has the same complexity order. Similarly, the operation of $\mathbf{X}_i^T \mathbf{G}$ has the same complexity. For solving the \mathbf{G} -subproblem, the major operation is the SVD of \mathbf{P} . This step is also not computationally heavy – what we ask for is an economy-size SVD of a very thin matrix (of size $L \times K$, $L \gg K$). This has a complexity order of $\mathcal{O}(LK^2)$ flops [26], which is light.

In terms of memory, all the terms involved (i.e., $\mathbf{Q}_i, \mathbf{G}_i^{(r)}, \mathbf{X}_i \mathbf{Q}_i, \mathbf{X}_i^T \mathbf{X}_i \mathbf{Q}_i$ and $\mathbf{X}_i^T \mathbf{G}_i^{(r)}$) only require $\mathcal{O}(LK)$ memory or less, but the eigen-decomposition-based solution needs $\mathcal{O}(M_i^2)$ and $\mathcal{O}(L^2)$ memory to store $(\mathbf{X}_i^T \mathbf{X}_i)^{-1}$ and \mathbf{M} , respectively. Note that K is usually very small compared to L and M_i and can be controlled by the designer.

We should mention that the updates of all the \mathbf{Q}_i ’s can be done in parallel, which allows easy distributed implementation. Suppose that each node in a network stores a view \mathbf{X}_i and computes \mathbf{Q}_i locally. The amount of information exchange after each update of \mathbf{Q}_i with other nodes is small – to be precise, the only thing needs to be exchanged is $\mathbf{X}_i \mathbf{Q}_i$, which has the size of $L \times K$.

IV. CONVERGENCE PROPERTIES

In this section, we study convergence properties of Algorithm 1. Note that the algorithm alternates between convex and nonconvex set-constrained subproblems and the subproblems may or may not be solved to optimality. Existing convergence analysis for exact and inexact block coordinate descent such as those in [29]–[32] results do not cover the all the considered cases, and thus the convergence properties are not obvious. For the purpose of discussion, we first define a critical point, or, a KKT point, of Problem (5). A KKT point (\mathbf{G}, \mathbf{Q}) satisfies the following optimality conditions:

$$\begin{cases} \mathbf{0} \in \nabla_{\mathbf{Q}} f(\mathbf{Q}, \mathbf{G}) + \partial_{\mathbf{Q}} g(\mathbf{Q}) \\ \mathbf{0} = \nabla_{\mathbf{G}} f(\mathbf{Q}, \mathbf{G}) + \mathbf{G}\mathbf{\Lambda}, \\ \mathbf{G}^T \mathbf{G} = \mathbf{I}, \end{cases} \quad (12)$$

where $f(\mathbf{Q}, \mathbf{G}) = \frac{1}{2} \sum_{i=1}^I \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2$, $\mathbf{\Lambda}$ is a Lagrangian multiplier, $\mathbf{Q} = [\mathbf{Q}_1^T, \dots, \mathbf{Q}_I^T]^T$ is a collection of all the \mathbf{Q}_i 's, and

$$\nabla_{\mathbf{Q}} f(\mathbf{Q}, \mathbf{G}) = \begin{bmatrix} \nabla_{\mathbf{Q}_1} f(\mathbf{Q}, \mathbf{G}) \\ \vdots \\ \nabla_{\mathbf{Q}_I} f(\mathbf{Q}, \mathbf{G}) \end{bmatrix}, \quad \partial_{\mathbf{Q}} g(\mathbf{Q}) = \begin{bmatrix} \partial_{\mathbf{Q}_1} g_1(\mathbf{Q}_1) \\ \vdots \\ \partial_{\mathbf{Q}_I} g_I(\mathbf{Q}_I) \end{bmatrix}. \quad (13)$$

and $\partial_{\mathbf{Q}_i} g_i(\mathbf{Q}_i)$ denotes a subgradient of the (possibly) non-smooth function $g_i(\mathbf{Q}_i)$.

Using the above definition, we first show that

Proposition 1 Assume that $\alpha_i \leq 1/L_i$ for all i , where $L_i = \lambda_{\max}(\mathbf{X}_i^T \mathbf{X}_i)$ is the largest eigenvalue of $\mathbf{X}_i^T \mathbf{X}_i$. Also assume that $g_i(\cdot)$ is a closed convex function and that the proximal operator in Line 6 of Algorithm 1 can be solved to optimality, $T \geq 1$, and $\gamma \in (0, 1]$. Then,

- The objective value of Problem (3) is non-increasing. In addition, every limit point of the solution sequence $\{\mathbf{G}^{(r)}, \{\mathbf{Q}_i^{(r)}\}_{i=1}^I\}_{r=0,1,\dots}$ is a KKT point of Problem (3).
- (Global Convergence) Further assume that \mathbf{X}_i and $\mathbf{Q}_i^{(0)}$ for $i = 1, \dots, I$ are bounded and $\text{rank}(\mathbf{X}_i) = M_i$. Then, the whole solution sequence converges to the set \mathcal{K} that consists of all the KKT points, i.e., $\lim_{r \rightarrow \infty} d^{(r)}(\mathcal{K}) \rightarrow 0$, where $d^{(r)}(\mathcal{K}) = \min_{\mathbf{Y} \in \mathcal{K}} \|(\mathbf{G}^{(r)}, \{\mathbf{Q}_i^{(r)}\}) - \mathbf{Y}\|_F$.

The monotonicity of the cost value sequence is a desirable property for very large-size problems. It ensures that the algorithm makes progress in each iteration, which is especially meaningful when the algorithm starts from a good but sub-optimal initialization, e.g., the solution given by the method in [10]. Proposition 1 (a) also characterizes the limit points of the solution sequence. According to Theorem 1, even only one proximal gradient step is performed in each iteration r , every limit point of the algorithm is a KKT point of Problem (5). As we demonstrated in the proof, Algorithm 1 has a similar flavor of *block successive upper bound minimization* (BSUM) [30] and the *block prox-linear* (BPL) framework [31], [32]. However, BSUM does not cover nonconvex constraints such as $\mathbf{G}^T \mathbf{G} = \mathbf{I}$, and BPL does not cover the $\gamma = 1$ case where the \mathbf{G} -subproblem is optimally solved. Therefore, the proofs in [30]–[32] cannot be applied to show convergence

of Algorithm 1. The (b) part of Proposition 1 establishes the convergence of the whole solution sequence – which is a much stronger result. The assumptions, on the other hand, are also more restrictive, where the views all have full column rank, which was not assumed in the (a) part.

Proposition 1 assures that the algorithm will approach a KKT point asymptotically. In practice, it is also meaningful to estimate the number of iterations that is needed for the algorithm reaching a neighborhood of a KKT point. To this end, let us define the following potential function:

$$Z^{(r,r+1)} = \sum_{t=0}^{T-1} \left\| \nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r,t)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}} g(\mathbf{Q}^{(r,t+1)}) \right\|_F^2 + \left\| \nabla_{\mathbf{G}} f(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r)}) + \mathbf{\Lambda}^{(r+1)} \mathbf{G}^{(r+1)} \right\|_F^2. \quad (14)$$

One can see that $Z^{(r,r+1)}$ is a value that is determined by two consecutive outer iterates of the algorithm. $Z^{(r,r+1)}$ has the following property:

Lemma 1 $Z^{(r,r+1)} \rightarrow 0$ implies that $(\mathbf{Q}^{(r)}, \mathbf{G}^{(r)})$ approaches a KKT point.

The proof of Lemma 1 is in Appendix B. As a result, we can use the value of $Z^{(r,r+1)}$ to measure how close is the current iterate to a KKT point, thereby estimating the iteration complexity. Following this rationale, we show that

Theorem 1 (Iteration Complexity) Assume that $\alpha_i < 1/L_i$, $0 < \gamma < 1$ and $T \geq 1$. Let $\delta > 0$ and J be the number of iterations needed for that $Z^{(r,r+1)} \leq \delta$ holds for the first time. Then, there exists a constant v such that $\delta \leq v/J-1$; that is, the algorithm converges to a KKT point at least sublinearly.

The proof of Theorem 1 is relegated to Appendix B. The major message revealed by Theorem 1 is that to reduce the optimality gap (measured by the Z -function) between the current iterate and a KKT point to δ by Algorithm 1, the number of iterations is at most $\mathcal{O}(1/\delta)$. We should mention that the proof is based on worst-case analysis, and the speed in practice maybe much faster. One subtle point that is worth mentioning is that the analysis in Theorem B holds when $\gamma < 1$ – it does not apply to the case where the \mathbf{G} -subproblem is optimally solved. This reflects some interesting facts in alternating optimization – when subproblems are handled in a more conservative way and using a controlled step size, convergence properties are in general nicer.

Proposition 1 and Theorem 1 characterize convergence properties of the proposed algorithm with a general regularization term $g_i(\cdot)$. It is also interesting to consider the special cases where $g_i(\cdot) = 0$ and $g_i(\cdot) = \mu_i \|\cdot\|_F^2$ – which correspond to the original MAX-VAR formulation and its “diagonal loaded” version [10]. These two cases are *optimally solvable* via eigen-decomposition. Since we posit the problem as a non-convex optimization problem in (3) and use alternating optimization to tackle it for dealing with large-scale cases, a natural question to ask is whether we have traded optimality for scalability. The answer is negative, and the proof is rooted in classic matrix computation theory. To explain, let us denote

$U_1 = U_M(:, 1 : K)$ and $U_2 = U_M(:, K + 1 : L)$ as the K principal eigenvectors of M and the eigenvectors spanning its orthogonal complement, respectively. Recall that our ultimate goal is to find G that is a basis of the range space of $U_M(:, 1 : K)$, denoted by $\mathcal{R}(U_M(:, 1 : K)) = \mathcal{R}_K(M)$. We show that

Lemma 2 (*Global Optimality*) *Denote the eigenvalues of $M \in \mathbb{R}^{L \times L}$ by $\lambda_1, \dots, \lambda_L$ in descending order. Consider $g_i(\cdot) = 0$ and $g_i(\cdot) = \mu_i \|\cdot\|_F^2$. Choose $\gamma = 1$. Under the above settings, assume that each subproblem in (6) is solved to optimality, $\text{rank}(X_i) = M_i$, $\lambda_K > \lambda_{K+1}$, and $\mathcal{R}(G^{(0)})$ is not orthogonal to any component in $\mathcal{R}_K(M)$, i.e.,*

$$\cos(\theta) = \min_{u \in \mathcal{R}_K(M), v \in \mathcal{R}(G^{(0)})} |u^T v| / (\|u\|_2 \|v\|_2) > 0.$$

Then, Algorithm (1) solves Problem (3) optimally when $r \rightarrow \infty$. In addition, $\mathcal{R}(G^{(r)})$ approaches $\mathcal{R}_K(M)$ linearly, i.e.,

$$\text{dist}(\mathcal{R}(G^{(r)}), \mathcal{R}_K(M)) \leq \tan(\theta) (\lambda_{K+1}/\lambda_K)^r,$$

where $\tan(\theta) = \|(U_2^T G^{(0)} (U_1^T G^{(0)})^{-1})\|_2$ and $\text{dist}(\mathcal{R}(G^{(r)}), \mathcal{R}_K(M)) = \|U_2^T G^{(r)}\|_2$.

The proof of Lemma 2 can be found in the appendices. The insight of the proof is as follows: The inner iterations (w.r.t. t) implicitly construct a term in the form of $M G^{(r)} \in \mathbb{R}^{L \times K}$ at iteration r if the Q -subproblem is solved to optimality, while the outer loop (w.r.t. r) amounts to a variation of the *Orthogonal Iteration* (OI) algorithm [26], which is a generalization of the power iteration that can estimate the K leading eigenvectors. As such, the outer loop inherits all the properties of the OI algorithm.

The result in Lemma 2 is clearly encouraging: For the MAX-VAR GCCA problem, we can maintain global optimality using low-complexity iterations without memory explosion. In addition, the outer iterations (indexed by r) converges linearly to a global optimal solution – which is rather favorable in terms of convergence rate. On the other hand, the result therein is hardly practical – it holds only if the Q -subproblem is solved to optimality which is usually not easy to do when the problem size is large. A natural question is that, if the Q -subproblem is solved inexactly, e.g., using a small number of iterations of T , can we expect similar convergence properties? To address this, we further show that

Theorem 2 *Under the same assumptions as those in Lemma 2 except that each subproblem in (6) is solved to an accuracy ϵ , i.e., $\|Q_i^{(r+1)} - \bar{Q}_i^{(r+1)}\|_2 \leq \epsilon$, where $\bar{Q}_i^{(r+1)} = (X_i^T X_i)^{-1} X_i^T G^{(r)}$. Then, after r iterations, we have*

$$\begin{aligned} & \text{dist}(\mathcal{R}(G^{(r)}), \mathcal{R}_K(M)) \\ & \leq \tan(\theta) (\lambda_{K+1}/\lambda_K)^r + C \left(\sum_{i=1}^I \lambda_{\max}(X_i) \epsilon \right), \end{aligned} \quad (15)$$

where C is a constant.

Theorem 2 makes much sense for inexact alternating optimization: in practice, solving the subproblem in (6) is not an

easy task, and one may want to stop early (e.g., using a small T). Theorem 2 ensures that if T is large enough to obtain a “good enough” approximation of the solution of Problem (6), the algorithm still converges *linearly* to the desired solution up to some accuracy loss.

V. NUMERICAL RESULTS

In this section, we use synthetic data and real experiments to showcase the effectiveness of the proposed algorithm. The experiments are carried out using MATLAB on a Linux server with 128GB RAM and 2.0GHz CPU cores.

A. Simulations

We generate the synthetic data in the following way: First, we let $Z \in \mathbb{R}^{L \times N}$ be a common latent factor of different views, where the entries of Z are drawn from the zero-mean i.i.d. Gaussian distribution and $L \geq N$. Then, a ‘mixing matrix’ $A_i \in \mathbb{R}^{N \times M}$ is multiplied to Z , resulting in $Y_i = Z A_i$. Finally, we add noise so that $X_i = Y_i + \sigma N_i$. Here, A_i and N_i are generated in the same way as Z . We apply the algorithm with diagonal loading (i.e., $g_i(\cdot) = \mu_i \|\cdot\|_F^2$) and let $\mu_i = 0.1$. For the synthetic data simulations, we employ the optimal solution that is based on eigen-decomposition and the multiview latent semantic analysis (MVLSA) algorithm that was proposed in [10] as baselines. Recall that MVLSA approximates the views using several leading singular values and vectors.

In Fig. 1(a), we use a small problem instance to examine the convergence properties of the proposed algorithm. Specifically, we let $(L, M, N, I) = (500, 25, 20, 3)$. For such a problem size, the optimal approach via eign-decomposition can be applied. We set $\sigma = 0.1$ in this case, let $P = 8$ for the MVLSA algorithm, and ask for $K = 5$ canonical components. The results are averaged over 50 random trials, where Z , $\{A_i\}$, $\{N_i\}$ are randomly generated in each trial. We test the proposed algorithm under different settings: We let $T = 1$, $T = 10$, and the GD run until the inner loop converges (denoted as ‘solved’ in the figures). We also initialize the algorithm with random initializations (denoted as ‘randn’) and warm starts (denoted as ‘warm’) – i.e., using the solutions of MVLSA as starting points. Some observations from Fig. 1(a) are in order. First, surprisingly, the proposed algorithm using various T ’s including $T = 1$ and random initialization can reach the global optimum, although Theorem 2 only covers the case where the Q_i -subproblem is solved to optimality. Second, by increasing T , the overall cost value decreases faster in terms of number of outer iterations – using $T = 10$ already gives very good speed of decreasing the cost value. Third, MVLSA cannot attain the global optimum, as expected. However, it provides good initialization: Using the warm start, the cost value comes close to the optimal value within 100 iterations in this case, even when $T = 1$ is employed. The cost value does not exhibit linear convergence, since in Theorem 2 the distance $\text{dist}(\mathcal{R}(G^{(r)}), \mathcal{R}_K(M))$ is defined by the matrix 2-norm (not Frobenius norm as in the cost function). In the Appendix, we show how $\text{dist}(\mathcal{R}(G^{(r)}), \mathcal{R}_K(M))$ evolves along the iterations, and the rate is indeed linear.

Fig. 1(b) shows the cost values against r when applying the algorithms to a relatively large-scale case. Here, we set $L = 100,000$, $M_i = 50,000$, $N = 1,000$, and $I = 3$. Each view are randomly generated and with 0.1% non-zero entries. For MVLSA, we let $P = 100$. The task is to seek $K = 10$ canonical components of the views and we stop the proposed algorithm when r reaches 100 and the result is averaged over 10 trials. Note that in this case, the optimal solution cannot be applied because of memory explosion. We see that the proposed algorithm with different settings converges to the same level within a few iterations. Our understanding is that the sparse views are well-conditioned and thus naturally lead to fast convergence of the GD part. We also show the runtimes of the algorithms when reaching the pointed cost levels in Fig. 1(b). Particularly, we see that using $T = 1$ exhibits very competitive runtime performance – Algorithm 1 with $T = 1$ reaches the pointed cost values using 6 seconds and 24 seconds with warm start and random initialization, respectively, while MVLSA uses 281 seconds to yield a worse solution.

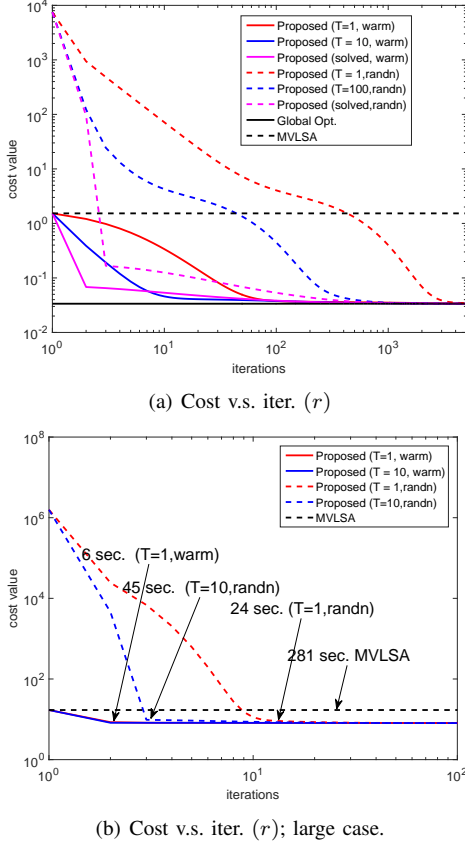


Fig. 1. Performance of the algorithms under various settings.

Fig. 2 shows the $\text{dist}(\mathcal{R}(\mathbf{G}^{(r)}), \mathcal{R}_K(\mathbf{M}))$'s obtained by the proposed algorithm under various settings against r . The settings are the same as those in Fig. 1(a) in the manuscript. We see that if the subproblem w.r.t. \mathbf{Q}_i is solved, the algorithm does give a linear convergence rate, as we stated in Theorem 2. Even we restrict the inner loop to $T = 10$ iterations, the rate is still empirically linear, but with a slightly worse slope.

The simulations with outliers go from here

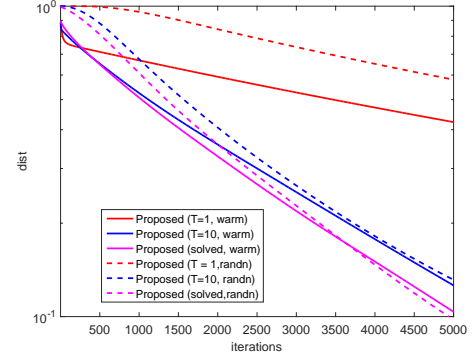


Fig. 2. Linear rate of subspace dist. v.s. iter. (r)

B. Real-Data Validation

We test the algorithms on a large-scale multilingual dataset. The views are extracted from a large word co-occurrence matrix, which is available at <https://sites.google.com/a/umn.edu/huang663/research>. The original data contains words of three languages, namely, English, Spanish, and French, and all the words are defined by co-occurrences with other words. We use the English words to form our first view, \mathbf{X}_1 , which contains $L = 183,034$ words and each word is defined by $M_i = 100,000$ features (co-occurrences). Note that \mathbf{X}_1 is sparse – only 1.21% of its entries are non-zeros. Using a dictionary, we pick out the translations of the English words contained in \mathbf{X}_1 in Spanish and French to form \mathbf{X}_2 and \mathbf{X}_3 , respectively. Note that many English words do not have a corresponding word in Spanish (or French). In such cases, we simply let $\mathbf{X}_i(\ell, :) = \mathbf{0}$ for $i = 2$ (or $i = 3$), resulting in sparser \mathbf{X}_2 and \mathbf{X}_3 .

Our objective is to find \mathbf{G} whose rows are the low-dimensional embeddings of the English words. To evaluate the output, we use the evaluation tool provided at wordvectors.org [33], which runs several word embedding tasks to evaluate a set of given embeddings. Simply speaking, the tasks compare the algorithm-learned embeddings with the judgment of humans and yield high scores if the embeddings are consistent with the humans. The scores are between zero and one, and a score equal to one means a perfect alignment between the learned result and human judgment. We use the result of MVLSA with $P = 640$ as benchmark. The result of applying SVD to \mathbf{X}_1 without considering different languages is also presented. We apply the proposed algorithm warm started by MVLSA and set $T = 1$. We run two versions of our algorithm. The first one uses $g_i(\cdot) = \|\cdot\|_F^2$ for $i = 1, 2, 3$. The second one uses $g_i(\cdot) = 0.05\|\cdot\|_{2,1}$ for $i = 2, 3$. The reason for adding ℓ_2/ℓ_1 mixed norm regularization to the French and Spanish views is twofold: First, the languages are effectively ‘fat matrices’ and thus need to use a column-selective regularizer. Second, the ℓ_2/ℓ_1 norm promotes row sparsity of \mathbf{Q}_i and thus performs feature selection on \mathbf{X}_2 and \mathbf{X}_3 – this physically means that we aim at selecting the most useful features from the other languages to help enhance English word embeddings.

Tables I and II show the word embedding results using $K = 50$ and $K = 100$, respectively. We see that using the

information from multiviews does help in improving the word embeddings: For $K = 50$ and $K = 100$, the multiview approaches perform better relative to SVD in 11 and 8 tasks out of 12 tasks. In addition, the proposed algorithm with the regularizer $g_i(\cdot) = \|\cdot\|_F^2$ gives similar or slightly better in average on both experiments compared to MVLSA. The proposed algorithm with the feature-selective regularizer ($g_i(\cdot) = \mu_i \|\cdot\|_{2,1}$) gives the best evaluation results on both experiments – this suggests that for large-scale multiview analysis, feature selection is much meaningful.

TABLE I
EVALUATION ON 12 WORD EMBEDDING TASKS; $K = 50$.

Task	Algorithm ($K = 50$)			
	svd	MVLSA	Proposed ($\ \cdot\ _F^2$)	Proposed ($\ \cdot\ _{2,1}$)
EN-WS-353-SIM	0.63	0.69	0.67	0.68
EN-MC-30	0.56	0.63	0.63	0.64
EN-MTurk-771	0.54	0.58	0.59	0.60
EN-MEN-TR-3k	0.67	0.66	0.67	0.68
EN-RG-65	0.51	0.53	0.55	0.58
EN-MTurk-287	0.65	0.64	0.65	0.64
EN-WS-353-REL	0.50	0.51	0.53	0.55
EN-VERB-143	0.21	0.22	0.21	0.21
EN-YP-130	0.36	0.39	0.38	0.41
EN-SIMLEX-999	0.31	0.42	0.41	0.39
EN-RW-STANFORD	0.39	0.43	0.43	0.43
EN-WS-353-ALL	0.56	0.59	0.59	0.60
average	0.49	0.52	0.53	0.54
median	0.53	0.56	0.57	0.59

TABLE II
EVALUATION ON 12 WORD EMBEDDING TASKS; $K = 100$.

Task	Algorithm ($K = 100$)			
	svd	MVLSA	Proposed ($\ \cdot\ _F^2$)	Proposed ($\ \cdot\ _{2,1}$)
EN-WS-353-SIM	0.68	0.72	0.71	0.72
EN-MC-30	0.73	0.68	0.72	0.74
EN-MTurk-771	0.59	0.60	0.60	0.61
EN-MEN-TR-3k	0.72	0.70	0.70	0.71
EN-RG-65	0.68	0.63	0.64	0.68
EN-MTurk-287	0.61	0.66	0.65	0.64
EN-WS-353-REL	0.57	0.54	0.55	0.56
EN-VERB-143	0.19	0.28	0.27	0.29
EN-YP-130	0.42	0.41	0.41	0.45
EN-SIMLEX-999	0.34	0.42	0.41	0.41
EN-RW-STANFORD	0.44	0.46	0.45	0.46
EN-WS-353-ALL	0.62	0.62	0.62	0.62
average	0.55	0.56	0.56	0.58
median	0.60	0.61	0.61	0.62

VI. CONCLUSION

In this work, we revisited the MAX-VAR GCCA problem with an eye towards scenarios involving large-scale and sparse data. The proposed AO-based approach is memory-efficient and has light per-iteration computational complexity if the views are sparse, and is thus suitable for dealing with big data. A thorough convergence analysis was presented, showing that the proposed algorithmic framework guarantees global optimality for the MAX-VAR GCCA problem when the subproblems in the AO framework are exactly solved in each iteration. In addition, when one subproblem is only inexactly solved, global convergence to a KKT point was also shown. Simulations and careful experiments with large-scale multi-lingual data showed that the performance of the proposed algorithm is promising in dealing with large and sparse multiview data.

APPENDIX A PROOF OF PROPOSITION 1

Before proving the Proposition, let us first simplify the notation. Recall that we have defined $\mathbf{Q} = [\mathbf{Q}_1^T, \dots, \mathbf{Q}_I^T]^T$ as a collection of \mathbf{Q}_i 's and we define $F(\mathbf{G}, \mathbf{Q}) = \sum_{i=1}^I \frac{1}{2} \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2 + \sum_{i=1}^I g_i(\mathbf{Q}_i)$, and the continuous differentiable part of the above as $f(\mathbf{G}, \mathbf{Q}) = \sum_{i=1}^I \frac{1}{2} \|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2$. Since the algorithm is essentially a two-block alternating optimization (i.e., \mathbf{Q}_i for all i are updated simultaneously), the above notation suffices to describe the updates. We also define

$$u_Q(\mathbf{Q}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) + \sum_{i=1}^I \langle \nabla_{\mathbf{Q}_i} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}), \mathbf{Q}_i - \hat{\mathbf{Q}}_i \rangle + \sum_{i=1}^I \frac{1}{2\alpha_i} \|\mathbf{Q}_i - \hat{\mathbf{Q}}_i\|_F^2 + \sum_{i=1}^I g_i(\mathbf{Q}_i);$$

i.e., $u_Q(\mathbf{Q}; \hat{\mathbf{G}}, \hat{\mathbf{Q}})$ is an approximation of $F(\mathbf{G}, \mathbf{Q})$ locally at the point $(\hat{\mathbf{G}}, \hat{\mathbf{Q}})$. We further define

$$\tilde{u}_Q(\mathbf{Q}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) + \sum_{i=1}^I \langle \nabla_{\mathbf{Q}_i} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}), \mathbf{Q}_i - \hat{\mathbf{Q}}_i \rangle + \sum_{i=1}^I \frac{1}{2\alpha_i} \|\mathbf{Q}_i - \hat{\mathbf{Q}}_i\|_F^2;$$

i.e., $\tilde{u}_Q(\mathbf{Q}; \hat{\mathbf{G}}, \hat{\mathbf{Q}})$ is an approximation of $f(\mathbf{G}, \mathbf{Q})$ locally at the point $(\hat{\mathbf{G}}, \hat{\mathbf{Q}})$. One can see that,

$$\nabla_{\mathbf{Q}_i} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) = \nabla_{\mathbf{Q}_i} \tilde{u}(\hat{\mathbf{Q}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}), \quad (16a)$$

$$\nabla_{\mathbf{Q}} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) + \partial_{\mathbf{Q}} g(\hat{\mathbf{Q}}) = \nabla_{\mathbf{Q}} \tilde{u}(\hat{\mathbf{Q}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) + \partial_{\mathbf{Q}} g(\hat{\mathbf{Q}}), \quad (16b)$$

where $\nabla_{\mathbf{Q}} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}})$, $\partial_{\mathbf{Q}} g(\mathbf{Q})$ and $\nabla_{\mathbf{Q}} \tilde{u}(\hat{\mathbf{Q}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}})$ follow the definitions in (13). Since $\nabla_{\mathbf{Q}_i} f(\mathbf{G}, \mathbf{Q})$ is L_i -Lipschitz continuous w.r.t. \mathbf{Q}_i and $\alpha_i \leq 1/L_i$ for all i , we have the following holds:

$$u_Q(\mathbf{Q}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) \geq F(\hat{\mathbf{G}}, \mathbf{Q}), \quad \forall \mathbf{Q}, \quad (17)$$

where the equality holds if and only if $\mathbf{Q}_i = \hat{\mathbf{Q}}_i$ for all i , i.e.,

$$u_Q(\hat{\mathbf{Q}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = F(\hat{\mathbf{G}}, \hat{\mathbf{Q}}). \quad (18)$$

Now, let us denote by $\mathbf{Q}^{(r,t)}$ (where $0 \leq t \leq T$) the solution of \mathbf{Q} after t gradient updates when $\mathbf{G}^{(r)}$ is fixed, where r is the iteration index of the outer loop. With the above notation, we have $\mathbf{Q}^{(r,0)} = \mathbf{Q}^{(r)}$ and $\mathbf{Q}^{(r,T)} = \mathbf{Q}^{(r+1)}$. Also, it can be seen that the update of \mathbf{Q}_i can be written as

$$\begin{aligned} \mathbf{Q}_i^{(r,t+1)} &= \text{prox}_g \left(\mathbf{Q}_i^{(r,t)} - \alpha_i \nabla_{\mathbf{Q}_i} f(\mathbf{G}^{(r)}, \mathbf{Q}^{(r,t)}) \right) \\ &= \arg \min_{\mathbf{Q}_i} u_Q(\mathbf{Q}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,t)}). \end{aligned} \quad (19)$$

Similarly, we define

$$u_G(\mathbf{G}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) + \langle \nabla_{\mathbf{G}} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}), \mathbf{G} - \hat{\mathbf{G}} \rangle + \frac{I}{2\gamma} \|\mathbf{G} - \hat{\mathbf{G}}\|_F^2 + \sum_{i=1}^I g_i(\hat{\mathbf{Q}}_i),$$

where the last term is a constant if \mathbf{Q} is fixed. We also have the following holds:

$$\nabla_{\mathbf{G}} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) = \nabla_{\mathbf{G}} u(\hat{\mathbf{G}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = I \cdot \hat{\mathbf{G}} - \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i. \quad (20)$$

The update rule of \mathbf{G} in Algorithm 1 can be re-expressed as follows:

$$\begin{aligned} \mathbf{G} &\in \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \left\| \mathbf{G} - \left((1-\gamma)\hat{\mathbf{G}} + \gamma \sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i / I \right) \right\|_F^2 \\ \Leftrightarrow \mathbf{G} &\in \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \left\| \mathbf{G} - \left(\hat{\mathbf{G}} - (\gamma/I) \nabla_{\mathbf{G}} f(\hat{\mathbf{G}}, \hat{\mathbf{Q}}) \right) \right\|_F^2 \\ \Leftrightarrow \mathbf{G} &\in \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} u_G(\mathbf{G}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) \end{aligned}$$

Since $\nabla_{\mathbf{G}} f(\mathbf{G}, \mathbf{Q})$ is I -Lipschitz continuous w.r.t. \mathbf{G} and $\gamma \leq 1$, we have

$$u_G(\mathbf{G}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) \geq F(\mathbf{G}, \hat{\mathbf{Q}}), \quad \forall \mathbf{G}, \quad (21)$$

and

$$u_G(\hat{\mathbf{G}}; \hat{\mathbf{G}}, \hat{\mathbf{Q}}) = F(\hat{\mathbf{G}}, \hat{\mathbf{Q}}). \quad (22)$$

Hence, Algorithm 1 boils down to

$$\mathbf{Q}_i^{(r,t+1)} = \arg \min_{\mathbf{Q}_i} u_Q(\mathbf{Q}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,t)}), \quad t = 0, \dots, T-1 \quad (23a)$$

$$\mathbf{G}^{(r+1)} \in \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} u_G(\mathbf{G}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,T)}). \quad (23b)$$

When $\gamma = 1$, (23b) amounts to SVD of $\sum_{i=1}^I \mathbf{X}_i \mathbf{Q}_i / I$ and the \mathbf{G} -subproblem is optimally solved.

Note that the following holds:

$$F(\mathbf{G}^{(r)}, \mathbf{Q}^{(r)}) = u_Q(\mathbf{Q}^{(r)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r)}) \quad (24a)$$

$$\geq u_Q(\mathbf{Q}^{(r,1)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,0)}) \quad (24b)$$

$$\geq F(\mathbf{G}^{(r)}, \mathbf{Q}^{(r,1)}) \quad (24c)$$

$$= u_Q(\mathbf{Q}^{(r,1)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,1)}) \quad (24d)$$

$$\geq u_Q(\mathbf{Q}^{(r,2)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,1)}) \quad (24e)$$

$$\vdots \quad (24f)$$

$$\geq u_Q(\mathbf{Q}^{(r+1)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r,T-1)}) \quad (24g)$$

$$\geq F(\mathbf{G}^{(r)}, \mathbf{Q}^{(r+1)}) \quad (24h)$$

$$= u_G(\mathbf{G}^{(r)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r+1)}) \quad (24i)$$

$$\geq u_G(\mathbf{G}^{(r+1)}; \mathbf{G}^{(r)}, \mathbf{Q}^{(r+1)}) \quad (24j)$$

$$\geq F(\mathbf{G}^{(r+1)}, \mathbf{Q}^{(r+1)}), \quad (24k)$$

where (24a) holds because of (18), (24b)-(24g) hold by invoking the update rule (23a) and the properties in (17) and (18),

(24h) holds because of (17), (24i) holds due to (22), (24j) is due to the fact that (23b) is optimally solved, and (24k) holds because of (21).

Next, we show that every limit point is a KKT point. Assume that there exists a convergent subsequence of $\{\mathbf{G}^{(r)}, \mathbf{Q}^{(r)}\}_{r=0,1,\dots}$, whose limit point is $(\mathbf{G}^*, \mathbf{Q}^*)$ and the subsequence is indexed by $\{r_j\}_{j=1,\dots,\infty}$. We have the following chain of inequalities:

$$u_Q(\mathbf{Q}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j)}) \geq u_Q(\mathbf{Q}^{(r_j,1)}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j)}) \quad (25a)$$

$$\geq u_Q(\mathbf{Q}^{(r_j,T)}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j,T-1)}) \quad (25b)$$

$$\geq F(\mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j+1)}) \quad (25c)$$

$$\geq F(\mathbf{G}^{(r_j+1)}, \mathbf{Q}^{(r_j+1)}) \quad (25d)$$

$$\geq F(\mathbf{G}^{(r_{j+1})}, \mathbf{Q}^{(r_{j+1})}) \quad (25e)$$

$$= u_Q(\mathbf{Q}^{(r_{j+1})}; \mathbf{G}^{(r_{j+1})}, \mathbf{Q}^{(r_{j+1})}), \quad (25f)$$

where (25a) holds because of the update rule in (23a), (25b) holds by repeating the arguments in (24b)-(24g), (25d) follows (24k), and (25f) is again because of the way that we construct $u_Q(\mathbf{Q}; \mathbf{G}^{(r_{j+1})}, \mathbf{Q}^{(r_{j+1})})$.

Taking $j \rightarrow \infty$, and by continuity of $u_Q(\cdot)$, we have

$$u_Q(\mathbf{Q}; \mathbf{G}^*, \mathbf{Q}^*) \geq u_Q(\mathbf{Q}^*; \mathbf{G}^*, \mathbf{Q}^*), \quad (26)$$

i.e., \mathbf{Q}^* is a minimum of $u_Q(\mathbf{Q}; \mathbf{G}^*, \mathbf{Q}^*)$. Consequently, \mathbf{Q}^* satisfies the conditional KKT conditions, i.e.,

$$\mathbf{0} \in \nabla_{\mathbf{Q}} \tilde{u}_Q(\mathbf{Q}^*; \mathbf{G}^*, \mathbf{Q}^*) + \partial_{\mathbf{Q}} g(\mathbf{Q}^*),$$

which also means that the following holds:

$$\mathbf{0} \in \nabla_{\mathbf{Q}_i} f(\mathbf{G}^*, \mathbf{Q}^*) + \partial_{\mathbf{Q}_i} g(\mathbf{Q}^*), \quad (27)$$

following (16).

We now show that $\mathbf{Q}^{(r_j,t)}$ for $t = 1, \dots, T$ also converges to \mathbf{Q}^* . Indeed, we have

$$\begin{aligned} u_Q(\mathbf{Q}^{(r_{j+1})}; \mathbf{G}^{(r_{j+1})}, \mathbf{Q}^{(r_{j+1})}) &\leq u_Q(\mathbf{Q}^{(r_j,1)}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j)}) \\ &\leq u_Q(\mathbf{Q}^{(r_j)}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j)}), \end{aligned}$$

where the first inequality was derived from (25). Taking $j \rightarrow \infty$, we see that $u_Q(\mathbf{Q}^*; \mathbf{G}^*, \mathbf{Q}^*) \leq u_Q(\mathbf{Q}^{(r_j,1)}; \mathbf{G}^*, \mathbf{Q}^*) \leq u_Q(\mathbf{Q}^*; \mathbf{G}^*, \mathbf{Q}^*)$, which implies that $u_Q(\mathbf{Q}^{(r_j,1)}; \mathbf{G}^*, \mathbf{Q}^*) = u_Q(\mathbf{Q}^*; \mathbf{G}^*, \mathbf{Q}^*) \leq u_Q(\mathbf{Q}; \mathbf{G}^*, \mathbf{Q}^*)$. On the other hand, the problem in (23a) has a unique minimizer when $g_i(\cdot)$ is a convex closed function [28], which means that $\mathbf{Q}^{(r_j,1)} \rightarrow \mathbf{Q}^*$. By recursion, we can show that $\mathbf{Q}^{(r_j,t)}$ for $t = 1, \dots, T$ (where T is finite) also converges to \mathbf{Q}^* using the same argument. Consequently, we have

$$\mathbf{Q}^{(r_j,T)} = \mathbf{Q}^{(r_{j+1})} \rightarrow \mathbf{Q}^*. \quad (28)$$

Now, we repeat the proof in (25) to \mathbf{G} :

$$\begin{aligned} u_G(\mathbf{G}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j+1)}) &\geq u_G(\mathbf{G}^{(r_j+1)}; \mathbf{G}^{(r_j)}, \mathbf{Q}^{(r_j+1)}) \\ &\geq F(\mathbf{G}^{(r_j+1)}, \mathbf{Q}^{(r_j+1)}) \\ &\geq F(\mathbf{G}^{(r_j+1)}, \mathbf{Q}^{(r_j+1)}) \\ &= u_G(\mathbf{G}^{(r_j+1)}; \mathbf{G}^{(r_j+1)}, \mathbf{Q}^{(r_j+1)}), \end{aligned}$$

Taking $j \rightarrow \infty$ and invoking (28), we have

$$u_G(\mathbf{G}; \mathbf{G}^*, \mathbf{Q}^*) \geq u_G(\mathbf{G}^*; \mathbf{G}^*, \mathbf{Q}^*), \quad \forall \mathbf{G}^T \mathbf{G} = \mathbf{I}.$$

The above means that \mathbf{G}^* satisfies the partial conditional KKT conditions w.r.t. \mathbf{G} . Combining with (27), we see that $(\mathbf{G}^*, \mathbf{Q}^*)$ is a KKT point of the original problem.

Now, we show the b) part. First, we show that \mathbf{Q}_i remains in a bounded set (the variable \mathbf{G} is always bounded since we keep it feasible in each iteration). Since the objective value is non-increasing (cf. Proposition 1), if we denote the initial objective value as V , then $F(\mathbf{G}^{(r)}, \mathbf{Q}^{(r)}) \leq V$ holds in all subsequent iterations. Note that when $\mathbf{X}_i^{(0)}$ and $\mathbf{Q}_i^{(0)}$ are bounded, V is also finite. In particular, we have $\|\mathbf{X}_i \mathbf{Q}_i - \mathbf{G}\|_F^2 + 2 \sum_{i=1}^I g_i(\mathbf{Q}_i) \leq 2V$ holds, which implies

$$\|\mathbf{X}_i \mathbf{Q}_i\|_F \leq \|\mathbf{G}\|_F + \sqrt{2V} \quad (30)$$

by the triangle inequality. The right-hand side of (30) is finite since both terms are bounded. Denote $(\|\mathbf{G}\|_F + \sqrt{2V})$ by V' . Then, we have

$$\begin{aligned} \|\mathbf{Q}_i\|_F &= \|(\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T \mathbf{X}_i \mathbf{Q}_i\|_F \\ &\leq \|(\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T\|_F \cdot \|\mathbf{X}_i \mathbf{Q}_i\|_F \\ &\leq V' \cdot \|(\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T\|_F. \end{aligned}$$

Now, by the assumption that $\text{rank}(\mathbf{X}_i) = M_i$, the term $\|(\mathbf{X}_i^T \mathbf{X}_i)^{-1} \mathbf{X}_i^T\|_F$ is bounded. This shows that $\|\mathbf{Q}_i\|_F$ is bounded. Hence, starting from a bounded $\mathbf{Q}_i^{(0)}$, the solution sequence $\{\{\mathbf{Q}_i^{(r)}\}, \mathbf{G}^{(r)}\}$ remains in a bounded set. Since the constraints of \mathbf{Q}_i , i.e., $\mathbb{R}^{M_i \times K}$ and \mathbf{G} are also closed sets, $\{\{\mathbf{Q}_i^{(r)}\}, \mathbf{G}^{(r)}\}$ remains in a compact set.

Now, let us denote \mathcal{K} as the set containing all the KKT points. Suppose the whole sequence does not converge to \mathcal{K} . Then, there exists a convergent subsequence indexed by $\{r_j\}$ such that $\lim_{j \rightarrow \infty} d(\mathcal{K}) \geq \gamma$ for some positive γ , where $d(\mathcal{K}) = \min_{\mathbf{Y} \in \mathcal{K}} \|(\mathbf{G}, \{\mathbf{Q}_i\}) - \mathbf{Y}\|$. Since the subsequence indexed by $\{r_j\}$ lies in a closed and bounded set as we have shown, this subsequence has a limit point. However, as we have shown in Theorem 1, every limit point of the solution sequence is a KKT point. This is a contradiction. Therefore, the whole sequence converges to a KKT point.

APPENDIX B

PROOF OF LEMMA 1 AND THEOREM 1

To keep the notation simple, we prove the theorem using the $T = 1$ case. The proof of the $T \geq 2$ can be obtained in a straightforward manner. We first show that $Z^{(r, r+1)} \rightarrow 0$ implies that a KKT point is reached. First, by the updating rule, we have

$$\begin{aligned} \mathbf{Q}^{(r, t)} &= \arg \min_{\mathbf{Q}} \left\langle \nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r, t-1)}, \mathbf{G}^{(r)}), \mathbf{Q} - \mathbf{Q}^{(r, t-1)} \right\rangle \\ &\quad + \sum_{i=1}^I g_i(\mathbf{Q}_i) + \sum_{i=1}^I \frac{1}{2\alpha_i} \|\mathbf{Q}_i - \mathbf{Q}_i^{(r, t-1)}\|_F^2. \end{aligned} \quad (31)$$

Therefore, there exists a $\partial_{\mathbf{Q}} g(\mathbf{Q}^{(r, t)})$, $\mathbf{Q}^{(r, t)}$ satisfies the following optimality conditions:

$$\mathbf{0} = \nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r, t-1)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}} g(\mathbf{Q}^{(r, t)}) + \mathbf{D}(\mathbf{Q}^{(r, t)} - \mathbf{Q}^{(r, t-1)}),$$

where $\mathbf{D} = \text{Diag}\left(\frac{1}{\alpha_1} \mathbf{1}_{M_1}^T, \dots, \frac{1}{\alpha_I} \mathbf{1}_{M_I}^T\right)$. Summing over the above through $t = 0$ to $t = T - 1$ we obtain

$$\begin{aligned} &\sum_{t=0}^{T-1} \left(\nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r, t)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}} g(\mathbf{Q}^{(r, t+1)}) \right) \\ &= - \left(\mathbf{D}(\mathbf{Q}^{(r)} - \mathbf{Q}^{(r+1)}) \right). \end{aligned} \quad (32)$$

Consequently, we have

$$\begin{aligned} &\sum_{t=0}^{T-1} \left\| \nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r, t)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}} g(\mathbf{Q}^{(r, t+1)}) \right\|_F^2 \rightarrow 0 \\ &\Rightarrow \mathbf{Q}_i^{(r, t)} - \mathbf{Q}_i^{(r, t+1)} \rightarrow \mathbf{0}, \quad \forall t = 0, \dots, T-1 \\ &\Rightarrow \mathbf{Q}_i^{(r)} - \mathbf{Q}_i^{(r+1)} \rightarrow \mathbf{0}, \end{aligned}$$

which holds since T is finite. The above means that $\mathbf{0} \in \nabla_{\mathbf{Q}} f(\mathbf{Q}^{(r)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}} g(\mathbf{Q}^{(r)})$ is satisfied when $Z^{r, r+1} \rightarrow 0$.

Recall that the update rule of \mathbf{G} is equivalent to solving

$$\min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \left\| \mathbf{G} - \left(\mathbf{G}^{(r)} - \gamma/I \left(\nabla_{\mathbf{G}} f(\mathbf{G}^{(r)}, \mathbf{Q}^{(r+1)}) \right) \right) \right\|_F^2. \quad (33)$$

Therefore, following the argument in (31), we have

$$\begin{aligned} \mathbf{G}^{(r+1)} &\in \arg \min_{\mathbf{G}^T \mathbf{G} = \mathbf{I}} \left\langle \nabla_{\mathbf{G}} f(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r)}), \mathbf{G} - \mathbf{G}^{(r)} \right\rangle \\ &\quad + \frac{1}{2\tilde{\gamma}} \|\mathbf{G} - \mathbf{G}^{(r)}\|_F^2, \end{aligned} \quad (34)$$

where $\tilde{\gamma} = \gamma/I$ and

$$\begin{aligned} &\nabla_{\mathbf{G}} f(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r)}) + \frac{1}{\tilde{\gamma}} \left(\mathbf{G}^{(r+1)} - \mathbf{G}^{(r)} \right) \\ &\quad + \mathbf{G}^{(r+1)} \mathbf{\Lambda}^{(r+1)} = \mathbf{0} \\ &\Rightarrow \left\| \nabla_{\mathbf{G}} f(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r)}) + \mathbf{G}^{(r+1)} \mathbf{\Lambda}^{(r+1)} \right\|_F^2 \\ &= \frac{1}{\tilde{\gamma}^2} \left\| \mathbf{G}^{(r+1)} - \mathbf{G}^{(r)} \right\|_F^2. \end{aligned} \quad (35)$$

Combining (35) and (32), we have

$$Z^{(r, r+1)} = \frac{1}{\tilde{\gamma}^2} \left\| \mathbf{G}^{(r+1)} - \mathbf{G}^{(r)} \right\|_F^2 + \sum_{i=1}^I \frac{1}{\alpha_i^2} \left\| \mathbf{Q}_i^{(r+1)} - \mathbf{Q}_i^{(r)} \right\|_F^2. \quad (36)$$

We see that $Z^{(r, r+1)} \rightarrow 0$ implies that $(\mathbf{G}^{(r+1)}, \mathbf{Q}^{(r+1)}) \rightarrow (\mathbf{G}^{(r)}, \mathbf{Q}^{(r)})$ and that a KKT point is reached.

Second, we show that every iterate of \mathbf{Q} and \mathbf{G} gives sufficiently large decrease of the overall objective function.

Since $\nabla_{Q_i} f(Q, G)$ is L_i -Lipschitz continuous for all i , we have the following:

$$\begin{aligned} F(Q^{(r,t+1)}, G^{(r)}) &\leq f(Q^{(r,t)}, G^{(r)}) \\ &\quad + \langle \nabla_Q f(Q^{(r,t)}, G^{(r)}), Q - Q^{(r,t)} \rangle \\ &\quad + \sum_{i=1}^I g_i(Q_i) + \sum_{i=1}^I \frac{L_i}{2} \|Q_i - Q_i^{(r,t)}\|_F^2. \end{aligned} \quad (37)$$

Since $Q^{(r,t)}$ is a minimizer of Problem (31), we also have

$$\begin{aligned} \langle \nabla_Q f(Q^{(r,t)}, G^{(r)}), Q^{(r,t+1)} - Q^{(r,t)} \rangle &+ \sum_{i=1}^I g_i(Q_i^{(r+1)}) \\ &+ \sum_{i=1}^I \frac{1}{2\alpha_i} \|Q_i^{(r,t+1)} - Q_i^{(r,t)}\|_F^2 \leq \sum_{i=1}^I g_i(Q_i^{(r,t+1)}). \end{aligned} \quad (38)$$

Combining (37) and (38), we have

$$\begin{aligned} F(Q^{(r,t+1)}, G^{(r)}) - F(Q^{(r,t)}, G^{(r)}) \\ \leq - \sum_{i=1}^I \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \|Q_i^{(r,t+1)} - Q_i^{(r,t)}\|_F^2. \end{aligned} \quad (39)$$

Summing up the above over $t = 0, \dots, T-1$,

$$\begin{aligned} F(Q^{(r)}, G^{(r)}) - F(Q^{(r+1)}, G^{(r)}) \\ \geq \sum_{t=0}^{T-1} \sum_{i=1}^I \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \|Q_i^{(r,t+1)} - Q_i^{(r,t)}\|_F^2. \end{aligned} \quad (40)$$

By the same derivation, we have

$$\begin{aligned} F(Q^{(r+1)}, G^{(r+1)}) - F(Q^{(r+1)}, G^{(r)}) \\ \leq - \left(\frac{1}{2\tilde{\gamma}} - \frac{1}{2} \right) \|G^{(r+1)} - G^{(r)}\|_F^2, \quad \forall G^T G = I. \end{aligned} \quad (41)$$

Combining (40) and (41), we have

$$\begin{aligned} F(Q^{(r)}, G^{(r)}) - F(Q^{(r+1)}, G^{(r+1)}) \\ \geq \left(\frac{1}{2\tilde{\gamma}} - \frac{1}{2} \right) \|G^{(r+1)} - G^{(r)}\|_F^2 \\ + \sum_{t=0}^{T-1} \sum_{i=1}^I \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \|Q_i^{(r,t+1)} - Q_i^{(r,t)}\|_F^2. \end{aligned} \quad (42)$$

Summing up $F(Q^{(r)}, G^{(r)})$ over $r = 0, 1, \dots, J-1$, we have Eq. (43), where $\alpha_{\min} = \min\{\alpha_1, \dots, \alpha_I\}$ and

$$c = \min \left\{ \left(\frac{1}{2\tilde{\gamma}} - \frac{1}{2} \right) \tilde{\gamma}^2, \left\{ \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \alpha_i^2 \right\}_{i=1, \dots, I} \right\}.$$

By the definition of J , we have

$$\begin{aligned} \frac{F(Q^{(0)}, G^{(0)}) - F(Q^{(J)}, G^{(J)})}{J-1} &\geq \frac{\sum_{r=0}^{J-1} c Z^{(r,r+1)}}{J-1} \geq c \cdot \epsilon \\ \Rightarrow \epsilon &\leq \frac{1}{c} \frac{F(Q^{(0)}, G^{(0)}) - \bar{F}}{J-1} \\ \Rightarrow \epsilon &\leq \frac{v}{J-1}, \end{aligned}$$

where \bar{F} is the lower bound of the cost function and

$$v = \frac{F(Q^{(0)}, G^{(0)}) - \bar{F}}{c}.$$

This completes the proof.

APPENDIX C PROOF OF LEMMA 2

If the subproblem w.r.t. Q_i is solved at outer iteration r , then, by the first-order optimality condition and the assumption that X_i has full column rank, we have

$$Q_i^{(r+1)} = (X_i^T X_i)^{-1} X_i^T G^{(r)}. \quad (44)$$

Therefore, the update w.r.t. G is simply

$$U_P^{(r)} \Sigma_P^{(r)} (V_P^{(r)})^T = \text{svd}(M G^{(r)}, \text{'econ'})} \quad (45a)$$

$$G^{(r+1)} = U_P^{(r)} (V_P^{(r)})^T. \quad (45b)$$

Since $V_P^{(r)} \in \mathbb{R}^{K \times K}$ is an orthonormal matrix, $G^{(r+1)}$ is an orthogonal basis of $M G^{(r)}$. In other words, there exists an invertible $\Theta^{(r+1)}$ such that

$$G^{(r+1)} \Theta^{(r+1)} = M G^{(r)}. \quad (46)$$

The update rule in (46), is essentially the orthogonal iteration algorithm in [26]. Invoking [26, Theorem 8.2.2], the proof is complete.

APPENDIX D PROOF OF THEOREM 2

The proof follows the insight of the alternating least squares-based generalized eigen-decomposition in [15] with modifications to accommodate the problem structure of MAX-VAR GCCA. At the r th iteration, ideally, we have $\tilde{Q}_i^{(r+1)} = (X_i^T X_i)^{-1} X_i^T G^{(r)}$ if the inner problem is solved to optimality. In practice, what we have is an inexact solution, i.e., $Q_i^{(r+1)} = (X_i^T X_i)^{-1} X_i^T G^{(r)} + W_i^{(r)}$, where we assume that the largest singular value of $W_i^{(r)}$ is bounded by ϵ , i.e., $\|W_i^{(r)}\|_2 \leq \epsilon$. Hence, we see that $\sum_{i=1}^I X_i^T Q_i^{(r+1)} = M G^{(r)} + \sum_{i=1}^I X_i W_i^{(r)}$. By the algorithm under $\gamma = 1$, we have

$$G^{(r+1)} = \left(M G^{(r)} + \sum_{i=1}^I X_i W_i^{(r)} \right) \Theta^{(r)},$$

where $\Theta^{(r)} \in \mathbb{R}^{K \times K}$ is a full-rank matrix since the solution via SVD is a change of bases. Let us denote U_1 and U_2 as orthogonal bases of $\mathcal{R}_K(M)$ and its orthogonal complement, respectively. Then, we have

$$\begin{bmatrix} U_1^T G^{(r+1)} \\ U_2^T G^{(r+1)} \end{bmatrix} = \begin{bmatrix} \Lambda_1 U_1^T G^{(r)} + U_1^T \sum_{i=1}^I X_i W_i^{(r)} \\ \Lambda_2 U_2^T G^{(r)} + U_2^T \sum_{i=1}^I X_i W_i^{(r)} \end{bmatrix} \Theta^{(r)}. \quad (47)$$

Now, we observe the equation on the top of the next page. Note that we can normalize the matrix $U_1^T \sum_{i=1}^I X_i W_i^{(r)}$ as follows

$$U_1^T \sum_{i=1}^I X_i W_i^{(r)} = \sigma \cdot \frac{U_1^T \sum_{i=1}^I X_i W_i^{(r)}}{\|U_1^T \sum_{i=1}^I X_i W_i^{(r)}\|_2} = \sigma \tilde{W}^{(r)}, \quad (49)$$

where σ is bounded by $\sigma \leq \sum_{i=1}^I \lambda_{\max}(X_i) \epsilon$ where ϵ denotes the 2-norm upper bound of $W_i^{(r)}$. Using the above notations, we come up with Eq. (50) on the next page, where the equality is obtained by the Taylor expansion.

$$\begin{aligned}
& F(\mathbf{Q}^{(r)}, \mathbf{G}^{(r)}) - F(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r+1)}) \\
& \geq \sum_{r=0}^{J-1} \left(\frac{1}{2\tilde{\gamma}} - \frac{1}{2} \right) \left\| \mathbf{G}^{(r+1)} - \mathbf{G}^{(r)} \right\|_F^2 + \sum_{r=0}^{J-1} \sum_{t=0}^{T-1} \sum_{i=1}^I \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \left\| \mathbf{Q}_i^{(r,t+1)} - \mathbf{Q}_i^{(r,t)} \right\|_F^2. \\
& = \sum_{r=0}^{J-1} \left(\frac{1}{2\tilde{\gamma}} - \frac{1}{2} \right) \tilde{\gamma}^2 \left\| \nabla_{\mathbf{G}} f(\mathbf{Q}^{(r+1)}, \mathbf{G}^{(r)}) + \mathbf{G}^{(r+1)} \boldsymbol{\Lambda}^{(r+1)} \right\|_F^2 \\
& + \sum_{r=0}^{J-1} \sum_{i=1}^I \sum_{t=0}^{T-1} \left(\frac{1}{2\alpha_i} - \frac{L_i}{2} \right) \alpha_i^2 \left\| \nabla_{\mathbf{Q}_i} f(\mathbf{Q}^{(r,t)}, \mathbf{G}^{(r)}) + \partial_{\mathbf{Q}_i} g_i(\mathbf{Q}^{(r,t+1)}) \right\|_F^2 \geq \sum_{r=0}^{J-1} c_Z^{(r,r+1)}, \quad (43)
\end{aligned}$$

$$\left\| \mathbf{U}_2^T \mathbf{G}^{(r+1)} \left(\mathbf{U}_1^T \mathbf{G}^{(r+1)} \right)^{-1} \right\|_2 = \left\| \left(\boldsymbol{\Lambda}_2 \mathbf{U}_2^T \mathbf{G}^{(r)} + \mathbf{U}_2^T \sum_{i=1}^I \mathbf{X}_i \mathbf{W}_i^{(r)} \right) \left(\boldsymbol{\Lambda}_2 \mathbf{U}_2^T \mathbf{G}^{(r)} + \mathbf{U}_2^T \sum_{i=1}^I \mathbf{X}_i \mathbf{W}_i^{(r)} \right)^{-1} \right\|_2. \quad (48)$$

$$\begin{aligned}
& \left\| \left(\boldsymbol{\Lambda}_2 \mathbf{U}_2^T \mathbf{G}^{(r)} + \mathbf{U}_2^T \sum_{i=1}^I \mathbf{X}_i \mathbf{W}_i^{(r)} \right) \left(\boldsymbol{\Lambda}_1 \mathbf{U}_1^T \mathbf{G}^{(r)} + \mathbf{U}_1^T \sum_{i=1}^I \mathbf{X}_i \mathbf{W}_i^{(r)} \right)^{-1} \right\|_2 \\
& = \left\| \left(\boldsymbol{\Lambda}_2 \mathbf{U}_2^T \mathbf{G}^{(r)} + \mathbf{U}_2^T \sum_{i=1}^I \mathbf{X}_i \mathbf{W}_i^{(r)} \right) \left(\left(\boldsymbol{\Lambda}_1 \mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} + \sigma \left(\boldsymbol{\Lambda}_1 \mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \tilde{\mathbf{W}}^{(r)} \left(\boldsymbol{\Lambda}_1 \mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} + \mathcal{O}(\sigma^2) \right) \right\|_2. \quad (50)
\end{aligned}$$

$$\left\| \mathbf{U}_2^T \mathbf{G}^{(r+1)} \left(\mathbf{U}_1^T \mathbf{G}^{(r+1)} \right)^{-1} \right\|_2 \leq \frac{\lambda_{K+1}}{\lambda_K} \left\| \left(\mathbf{U}_2^T \mathbf{G}^{(r)} \right) \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2 + \left(\sum_{i=1}^I \lambda_{\max}(\mathbf{X}_i) \epsilon \right) \mathcal{O} \left(\left\| \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2^2 \right). \quad (51)$$

Now, let us drop the second- and higher-order terms of σ which are sufficiently small and ‘absorb’ them in $\mathcal{O}(\left\| \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2^2)$. Consequently, we obtain Eq. (51). Now, we show that $\left\| \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2^2$ is bounded for all r . This can be seen by induction. For $r = 1$, we see that $\left\| \mathbf{U}_2^T \mathbf{G}^{(1)} \left(\mathbf{U}_1^T \mathbf{G}^{(1)} \right)^{-1} \right\|_2^2$ has to be bounded since we assumed that $\text{rank}(\mathbf{U}_1^T \mathbf{G}^{(0)}) = K$ and since (51) holds. Using the same argument, we see that for all $r \geq 1$, $\left(\mathbf{U}_1^T \mathbf{G}^{(1)} \right)^{-1}$ is bounded. Let us denote an upper bound as β , i.e.,

$$\left\| \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2 \leq \beta, \quad \forall r.$$

Then, we obtain Eq. (52). Using the above, we see Eq. (53) – which completes the proof.

REFERENCES

- [1] D. R. Hardoon, S. Szedmak, and J. Shawe-Taylor, “Canonical correlation analysis: An overview with application to learning methods,” *Neural computation*, vol. 16, no. 12, pp. 2639–2664, 2004.
- [2] Y.-O. Li, T. Adali, W. Wang, and V. D. Calhoun, “Joint blind source separation by multiset canonical correlation analysis,” *IEEE Trans. Signal Process.*, vol. 57, no. 10, pp. 3918–3929, 2009.
- [3] A. Bertrand and M. Moonen, “Distributed canonical correlation analysis in wireless sensor networks with application to distributed blind source separation,” *IEEE Trans. Signal Process.*, vol. 63, no. 18, pp. 4800–4813, 2015.
- [4] Q. Wu and K. M. Wong, “Un-music and un-cle: An application of generalized correlation analysis to the estimation of the direction of arrival of signals in unknown correlated noise,” *IEEE Trans. Signal Process.*, vol. 42, no. 9, pp. 2331–2343, 1994.
- [5] A. Dogandzic and A. Nehorai, “Finite-length mimo equalization using canonical correlation analysis,” *IEEE Trans. Signal Process.*, vol. 50, no. 4, pp. 984–989, 2002.
- [6] S. M. Kakade and D. P. Foster, “Multi-view regression via canonical correlation analysis,” in *Learning Theory*. Springer, 2007, pp. 82–96.
- [7] K. Chaudhuri, S. M. Kakade, K. Livescu, and K. Sridharan, “Multi-view clustering via canonical correlation analysis,” in *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009, pp. 129–136.
- [8] R. Arora and K. Livescu, “Multi-view learning with supervision for transformed bottleneck features,” in *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. IEEE, 2014, pp. 2499–2503.
- [9] H. Wang, T. Lee, C.-C. Leung, B. Ma, and H. Li, “Acoustic segment modeling with spectral clustering methods,” *Audio, Speech, and Language Processing, IEEE/ACM Transactions on*, vol. 23, no. 2, pp. 264–277, 2015.
- [10] P. Rastogi, B. Van Durme, and R. Arora, “Multiview LSA: Representation learning via generalized cca,” in *North American Chapter of the Association for Computational Linguistics (NAACL)*, 2015.
- [11] J. D. Carroll, “Generalization of canonical correlation analysis to three or more sets of variables,” in *Proceedings of the 76th annual convention of the American Psychological Association*, vol. 3, 1968, pp. 227–228.
- [12] J. R. Kettenring, “Canonical analysis of several sets of variables,” *Biometrika*, vol. 58, no. 3, pp. 433–451, 1971.
- [13] Z. Ma, Y. Lu, and D. Foster, “Finding linear structure in large datasets with scalable canonical correlation analysis,” *arXiv preprint arXiv:1506.08170*, 2015.
- [14] L. Sun, S. Ji, and J. Ye, “Canonical correlation analysis for multilabel classification: A least-squares formulation, extensions, and analysis,”

$$\left\| \mathbf{U}_2^T \mathbf{G}^{(r+1)} \left(\mathbf{U}_1^T \mathbf{G}^{(r+1)} \right)^{-1} \right\|_2 \leq \frac{\lambda_{K+1}}{\lambda_K} \left\| \left(\mathbf{U}_2^T \mathbf{G}^{(r)} \right) \left(\mathbf{U}_1^T \mathbf{G}^{(r)} \right)^{-1} \right\|_2 + \left(\sum_{i=1}^I \lambda_{\max}(\mathbf{X}_i) \epsilon \right) \mathcal{O}(\beta^2). \quad (52)$$

$$\left\| \mathbf{U}_2^T \mathbf{G}^{(r+1)} \left(\mathbf{U}_1^T \mathbf{G}^{(r+1)} \right)^{-1} \right\|_2 \leq \left(\frac{\lambda_{K+1}}{\lambda_K} \right)^r \left\| \left(\mathbf{U}_2^T \mathbf{G}^{(0)} \right) \left(\mathbf{U}_1^T \mathbf{G}^{(0)} \right)^{-1} \right\|_2 + \sum_{t=0}^{r-1} \left(\frac{\lambda_{K+1}}{\lambda_K} \right)^t \left(\sum_{i=1}^I \lambda_{\max}(\mathbf{X}_i) \epsilon \right) \mathcal{O}(\beta^2). \quad (53)$$

Pattern Analysis and Machine Intelligence, IEEE Transactions on, vol. 33, no. 1, pp. 194–200, 2011.

- [15] Y. Lu and D. P. Foster, “Large scale canonical correlation analysis with iterative least squares,” in *Advances in Neural Information Processing Systems*, 2014, pp. 91–99.
- [16] M. Van De Velden and T. H. A. Bijmolt, “Generalized canonical correlation analysis of matrices with missing rows: a simulation study,” *Psychometrika*, vol. 71, no. 2, pp. 323–331, 2006.
- [17] D. R. Hardoon and J. Shawe-Taylor, “Sparse canonical correlation analysis,” *Machine Learning*, vol. 83, no. 3, pp. 331–353, 2011.
- [18] D. M. Witten, R. Tibshirani, and T. Hastie, “A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis,” *Biostatistics*, p. kxp008, 2009.
- [19] X. Chen, H. Liu, and J. G. Carbonell, “Structured sparse canonical correlation analysis,” in *International Conference on Artificial Intelligence and Statistics*, 2012, pp. 199–207.
- [20] D. M. Witten and R. J. Tibshirani, “Extensions of sparse canonical correlation analysis with applications to genomic data,” *Statistical applications in genetics and molecular biology*, vol. 8, no. 1, pp. 1–27, 2009.
- [21] M. Faruqui and C. Dyer, “Improving vector space word representations using multilingual correlation.” Association for Computational Linguistics, 2014.
- [22] I. Rustandi, M. A. Just, and T. Mitchell, “Integrating multiple-study multiple-subject fmri datasets using canonical correlation analysis,” in *Proceedings of the MICCAI 2009 Workshop: Statistical modeling and detection issues in intra-and inter-subject functional MRI data analysis*, 2009.
- [23] T. M. Mitchell, S. V. Shinkareva, A. Carlson, K.-M. Chang, V. L. Malave, R. A. Mason, and M. A. Just, “Predicting human brain activity associated with the meanings of nouns,” *science*, vol. 320, no. 5880, pp. 1191–1195, 2008.
- [24] A. Tenenhaus and M. Tenenhaus, “Regularized generalized canonical correlation analysis,” *Psychometrika*, vol. 76, no. 2, pp. 257–284, 2011.
- [25] A. Tenenhaus, C. Philippe, V. Guillemot, K.-A. Le Cao, J. Grill, and V. Frouin, “Variable selection for generalized canonical correlation analysis,” *Biostatistics*, p. kxu001, 2014.
- [26] G. H. Golub and C. F. V. Loan., *Matrix Computations*. The Johns Hopkins University Press, 1996.
- [27] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends in Machine Learning*, pp. 1–122, 2011.
- [28] N. Parikh and S. Boyd, “Proximal algorithms,” *Foundations and Trends in optimization*, vol. 1, no. 3, pp. 123–231, 2013.
- [29] D. P. Bertsekas, *Nonlinear programming*. Athena Scientific, 1999.
- [30] M. Razaviyayn, M. Hong, and Z.-Q. Luo, “A unified convergence analysis of block successive minimization methods for nonsmooth optimization,” *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1126–1153, 2013.
- [31] Y. Xu and W. Yin, “A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion,” *SIAM Journal on imaging sciences*, vol. 6, no. 3, pp. 1758–1789, 2013.
- [32] —, “A globally convergent algorithm for nonconvex optimization based on block coordinate update,” *arXiv preprint arXiv:1410.1386*, 2014.
- [33] M. Faruqui and C. Dyer, “Community evaluation and exchange of word vectors at wordvectors.org,” in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstra-*

tions. Baltimore, USA: Association for Computational Linguistics, June 2014.