

Efficient Globally Convergent Stochastic Optimization for Canonical Correlation Analysis

Authored by:

Weiran Wang
Dan Garber
Dan Garber
Nati Srebro
Jialei Wang

Abstract

We study the stochastic optimization of canonical correlation analysis (CCA), whose objective is nonconvex and does not decouple over training samples. Although several stochastic gradient based optimization algorithms have been recently proposed to solve this problem, no global convergence guarantee was provided by any of them. Inspired by the alternating least squares/power iterations formulation of CCA, and the shift-and-invert preconditioning method for PCA, we propose two globally convergent meta-algorithms for CCA, both of which transform the original problem into sequences of least squares problems that need only be solved approximately. We instantiate the meta-algorithms with state-of-the-art SGD methods and obtain time complexities that significantly improve upon that of previous work. Experimental results demonstrate their superior performance.

1 Paper Body

Canonical correlation analysis (CCA, [1]) and its extensions are ubiquitous techniques in scientific research areas for revealing the common sources of variability in multiple views of the same phenomenon. In CCA, the training set consists of paired observations from two views, denoted $(x_1, y_1), \dots, (x_N, y_N)$, where N is the training set size, $x_i \in \mathbb{R}^{d_x}$ and $y_i \in \mathbb{R}^{d_y}$ for $i = 1, \dots, N$. We also denote the data matrices for each view by $X = [x_1, \dots, x_N] \in \mathbb{R}^{d_x \times N}$ and $Y = [y_1, \dots, y_N] \in \mathbb{R}^{d_y \times N}$, and $d := d_x + d_y$. The objective of CCA is to find linear projections of each view such that the correlation between the projections is maximized: $\max_{u, v}$

$$\begin{aligned} & u^T \Sigma_{xy} v \\ \text{s.t.} \end{aligned}$$

$$u^T \Sigma_{xx} u = v^T \Sigma_{yy} v = 1$$

where $\Sigma_{xy} = N^{-1} XY^T$ is the cross-covariance matrix, $\Sigma_{xx} = N^{-1} XX^T + \lambda_x I$ and $\Sigma_{yy} = N^{-1} YY^T + \lambda_y I$ are the auto-covariance matrices, and $(\lambda_x, \lambda_y) \geq 0$ are regularization parameters [2].

$$(1) \quad \max_{u, v} \frac{1}{N} XY^T uv^T +$$

We denote by (u^*, v^*) the global optimum of (1), which can be computed in closed-form. Define ρ^*

$$\rho^*$$

$$T := \Sigma_{xx}^2 \Sigma_{xy} \Sigma_{yy}^2 \int_0^1 R dx \int_0^1 dy,$$

$$(2)$$

and let (u_i, v_i) be the (unit-length) left and right singular vector pair associated with T 's largest singular value ρ_i . Then the optimal objective value, i.e., the canonical correlation between the ρ_1 views, is ρ^* , achieved by $(u^*, v^*) = (\rho_1^{-1/2} \Sigma_{xx}^{-1/2} u_1, \rho_1^{-1/2} \Sigma_{yy}^{-1/2} v_1)$. Note that

$$\rho_1$$

$$\rho_1 \rho_1 = k T k \leq \Sigma_{xx}^2 X$$

$\Sigma_{yy}^2 Y \leq 1$. Furthermore, we are guaranteed to have $\rho_i \leq 1$ if $(\lambda_x, \lambda_y) \geq 0$. The first two authors contributed equally. We assume that X and Y are centered at the origin for notational simplicity; if they are not, we can center them as a pre-processing operation. \square

$$2$$

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

Table 1: Time complexities of different algorithms for achieving ϵ -suboptimal solution (u, v) to

CCA, i.e., $\min_{u, v} (u^T \Sigma_{xx} u)^{-1/2}, (v^T \Sigma_{yy} v)^{-1/2} \geq 1 - \epsilon$. GD=gradient descent, AGD=accelerated GD, SVRG=stochastic variance reduced gradient, ASVRG=accelerated SVRG. Note ASVRG provides speedup over SVRG only when $\epsilon \leq N^{-1}$, and we show the dominant term in its complexity. Algorithm Least squares solver Time complexity

$$\frac{1}{\epsilon^2} \frac{1}{N} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ AppGrad [3] GD (local) } \frac{2}{\epsilon} \log \frac{1}{\epsilon}$$

$$\frac{1}{\epsilon} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ CCALin [6] AGD } \frac{1}{\epsilon^2} \frac{1}{N} \log \frac{1}{\epsilon} \frac{1}{\epsilon}$$

$$2$$

$$\frac{1}{\epsilon^2} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ This work: AGD } \frac{1}{\epsilon^2} \frac{1}{N} \log \frac{1}{\epsilon} \frac{1}{\epsilon}$$

Alternating least

$$\frac{1}{\epsilon^2} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ squares (ALS) SVRG } O\left(\frac{1}{\epsilon^2} \frac{d(N + \frac{1}{\epsilon})}{N}\right) \log \frac{1}{\epsilon} \frac{1}{\epsilon}$$

$$2$$

$$\frac{1}{\epsilon^2} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ ASVRG } \frac{1}{\epsilon^2} \frac{1}{N} \log \frac{1}{\epsilon} \frac{1}{\epsilon}$$

$$\frac{1}{\epsilon} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \text{ This work: AGD } \frac{1}{\epsilon^2} \frac{1}{N} \log \frac{1}{\epsilon} \frac{1}{\epsilon}$$

$\frac{1}{\epsilon}$ Shift-and-invert $\frac{1}{\epsilon} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \log \frac{1}{\epsilon} \frac{1}{\epsilon}$ SVRG preconditioning (SI)

$\frac{1}{\epsilon} \frac{1}{N} \frac{1}{\epsilon} dN = O\left(\frac{1}{\epsilon^2} \frac{dN}{N}\right) \log \frac{1}{\epsilon} \frac{1}{\epsilon}$ ASVRG $\frac{1}{\epsilon^2} \frac{1}{N} \log \frac{1}{\epsilon} \frac{1}{\epsilon}$ For large and high dimensional datasets, it is time and memory consuming to first explicitly form the matrix T (which requires eigen-decomposition of the covariance matrices) and then compute its singular value decomposition (SVD). For such datasets,

it is desirable to develop stochastic algorithms that have efficient updates, converge fast, and takes advantage of the input sparsity. There have been recent attempts to solve (1) based on stochastic gradient descent (SGD) methods [3, 4, 5], but none of these work provides rigorous convergence analysis for their stochastic CCA algorithms. The main contribution of this paper is the proposal of two globally convergent meta-algorithms for solving (1), namely, alternating least squares (ALS, Algorithm 2) and shift-and-invert preconditioning (SI, Algorithm 3), both of which transform the original problem (1) into sequences of least squares problems that need only be solved approximately. We instantiate the meta algorithms with state-of-the-art SGD methods and obtain efficient stochastic optimization algorithms for CCA. In order to measure the alignments between an approximate solution (u, v) and the optimum (u^*, v^*) , we assume that T has a positive singular value gap $\gamma := \lambda_1 - \lambda_2 \in (0, 1]$ so its top left and right singular vector pair is unique (up to a change of sign). Table 1 summarizes the time complexities of several algorithms for achieving ϵ -suboptimal alignment $\max(\|x\|_2, \|y\|_2) \leq \epsilon$ in terms, where $\epsilon = \min(\epsilon_{xx}, \epsilon_{yy})^{1/3}$ to hide poly-logarithmic dependencies (see problems solved in all cases). We use the notation $O(\cdot)$ Sec. 3.1.1 and Sec. 3.2.3 for the hidden factors). Each time complexity may be preferable in certain regime depending on the parameters of the problem. Notations We use $\lambda_i(A)$ to denote the i -th largest singular value of a matrix A , and use $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ to denote the largest and smallest singular values of A respectively.

2

Motivation: Alternating least squares

Our solution to (1) is inspired by the alternating least squares (ALS) formulation of CCA [7, Algorithm 5.2], as shown in Algorithm 1. Let the nonzero singular values of T be $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0$, where $r = \text{rank}(T) \leq \min(d_x, d_y)$, and the corresponding (unit-length) left and right singular vector pairs be $(a_1, b_1), \dots, (a_r, b_r)$, with $a_1 = \lambda_1^{-1/2}$ and $b_1 = \lambda_1^{-1/2}$. Define

$$U = T C = \frac{1}{\sqrt{\lambda_1}} R d^T d. \quad (3) \quad T^T \geq 0$$

For the ALS meta-algorithm, it's enough to consider a per-view conditioning. And when using AGD as the least squares solver, the time complexities depends on $\lambda_{\max}(\lambda_{xx})$ instead, which is less than $\max_i \lambda_{xi}^2$.

2

Algorithm 1 Alternating least squares for CCA. Input: Data matrices $X \in \mathbb{R}^{d_x \times N}$, $Y \in \mathbb{R}^{d_y \times N}$, regularization parameters $(\lambda_x, \lambda_y) \in \mathbb{R}_{>0}^2$. $\lambda_x \geq \lambda_y$. $\lambda_x \geq \lambda_y$. Initialize $u, v \in \mathbb{R}^{d_x}, \mathbb{R}^{d_y}$ $u \leftarrow \frac{1}{\sqrt{\lambda_x}}, v \leftarrow \frac{1}{\sqrt{\lambda_y}}$. For $t = 1, 2, \dots, T$ do $u \leftarrow \frac{1}{\sqrt{\lambda_x}} \frac{X^T Y v}{X^T X u + \lambda_x I}$, $v \leftarrow \frac{1}{\sqrt{\lambda_y}} \frac{Y^T X u}{Y^T Y v + \lambda_y I}$. Output: $(u^T, v^T) \approx (u^*, v^*)$ as $T \rightarrow \infty$.

on \mathbb{R}^n

?

?

$\lambda_{xx} \geq \lambda_{yy} \geq \lambda_{xy} \geq \lambda_{yx} \geq \lambda_{xx} \geq \lambda_{yy} \geq \lambda_{xy} \geq \lambda_{yx}$

n

exists no singular value gap, the top singular vector pair is not unique and it is no longer meaningful to measure the alignments. Nonetheless, it is possible to extend our proof to obtain sublinear convergence for the objective in this case. Observe that, besides the steps of normalization to unit length, the basic operation in each iteration 1 of Algorithm 1 is of the form $u_{xx} u_{xy} v_{t+1} = (N XX + \lambda I)^{-1} N XY v_t$, which is equivalent to solving the following regularized least squares (ridge regression) problem

$$\min_u \|X u - y\|^2 + \lambda \|u\|^2$$

where $X = X^T v_t$ and $y = y^T v_t$. In the next section, we show that, to maintain the convergence of ALS, it is unnecessary to solve the least squares problems exactly. This enables us to use state-of-the-art SGD methods for solving (6) to sufficient accuracy, and to obtain a globally convergent stochastic algorithm for CCA.

One can show that $\langle u_0, v_0 \rangle$ is bounded away from 0 with high probability using random initialization (u_0, v_0) .

3

Algorithm 2 The alternating least squares (ALS) meta-algorithm for CCA.
Input: Data matrices $X \in \mathbb{R}^{d \times N}$, $Y \in \mathbb{R}^{d \times N}$, regularization parameters (λ_x, λ_y) . Initialize u, v to random unit vectors. For $t = 1, 2, \dots, T$ do

1 $X^T X u + \lambda_x u = X^T Y v$

2 Solve for u with initialization u . $Y^T Y v + \lambda_y v = Y^T X u$

3 Solve for v with initialization v

4 $X^T X u + \lambda_x u = X^T Y v$, and output (u, v) satisfying $\|X u - Y v\|^2 + \lambda_x \|u\|^2 + \lambda_y \|v\|^2 \leq \epsilon$. approximate solution (u, v) is the approximate solution to CCA.

3.1

Our algorithms Algorithm I: Alternating least squares (ALS) with variance reduction

Our first algorithm consists of two nested loops. The outer loop runs inexact power iterations while the inner loop uses advanced stochastic optimization methods, e.g., stochastic variance reduced gradient (SVRG, [9]) to obtain approximate matrix-vector multiplications. A sketch of our algorithm is provided in Algorithm 2. We make the following observations from this algorithm. Connection to previous work At step t , if we optimize $f_t(u)$ and $g_t(v)$ crudely by a single batch step, we obtain the following update rule: gradient descent step from the initialization (u_t, v_t) , $u_{t+1} = (X^T X + \lambda_x I)^{-1} X^T Y v_t$, $v_{t+1} = (Y^T Y + \lambda_y I)^{-1} Y^T X u_t$ where λ is the stepsize (assuming $\lambda_x = \lambda_y = 0$). This coincides with the AppGrad algorithm of [3, Algorithm 3], for which only local convergence is shown. Since the objectives $f_t(u)$ and $g_t(v)$ decouple over training samples, it is convenient to apply SGD methods to

them. This observation motivated the stochastic CCA algorithms of [3, 4]. We note however, no global convergence guarantee was shown for these stochastic CCA algorithms, and the key to our convergent algorithm is to solve the least squares problems to sufficient accuracy. Warm-start Observe that for different t , the least squares problems $f_t(u)$ only differ in their targets as v_t changes over time. Since v_{t+1} is close to v_t (especially when near convergence), we may use u_t as initialization for minimizing $f_{t+1}(u)$ with an iterative algorithm. Normalization p At the end of each outer loop, Algorithm 2 implements exact normalization of the $u_t / \|u_t\| = N^{-1} \sum_{i=1}^N u_t^{(i)}$ to ensure the constraints, where $u_t^{(i)}$ form $u_t = \sum_{i=1}^N u_t^{(i)} X^{(i)} X^{(i)T} + \sum_{i=1}^N u_t^{(i)} X^{(i)}$. However, this does not introduce extra computation because we also compute this projection for the batch gradient used by SVRG (at the beginning of time step $t + 1$). In contrast, the stochastic algorithms of [3, 4] (possibly adaptively) estimate the covariance matrix from a minibatch of training samples and use the estimated covariance for normalization. This is because their algorithms perform normalizations after each update and thus need to avoid computing the projection of the entire training set frequently. But as a result, their inexact normalization steps introduce noise to the algorithms. Input sparsity For high dimensional sparse data (such as those used in natural language processing [10]), an advantage of gradient based methods over the closed-form solution is that the former takes into account the input sparsity. For sparse inputs, the time complexity of our algorithm depends on $\text{nnz}(X, Y)$, i.e., the total number of nonzeros in the inputs instead of dN . Canonical ridge When $(x, y) \neq 0$, $f_t(u)$ and $g_t(v)$ are guaranteed to be strongly convex due to the ℓ_2 regularizations, in which case SVRG converges linearly. It is therefore beneficial to use

small nonzero regularization for improved computational efficiency, especially for high dimensional datasets where inputs X and Y are approximately low-rank. Convergence By the analysis of inexact power iterations where the least squares problems are solved (or the matrix-vector multiplications are computed) only up to necessary accuracy, we provide the following theorem for the convergence of Algorithm 2 (see its proof in Appendix B). The key to our analysis is to bound the distances between the iterates of Algorithm 2 and that of Algorithm 1 at all time steps, and when the errors of the least squares problems are sufficiently small (at the level of ϵ^2), the iterates of the two algorithms have the same quality. Theorem 2 (Convergence of Algorithm 2). Fix $T \geq 2 \log \frac{1}{\epsilon^2}$, and set $\eta(T) \geq 2$

$\frac{2}{\eta(T)} \geq 2r \left(\frac{2}{\eta(T)} \right)$ in Algorithm 2. Then we have $u_T = v_T^T \eta(T) v_T = 1, T \geq 128 \frac{2}{\eta(T)}$

$\frac{2}{\eta(T)} \geq 2 \frac{1}{\eta(T)}$, and $u_T^T \min(u_T^T X X^T u_T), (v_T^T Y Y^T v_T)^T X Y^T v_T \geq 1$ (1 ≥ 2).

3.1.1 Stochastic optimization of regularized least squares We now discuss the inner loop of Algorithm 2, which approximately solves problems of the form (6). Owing to the finite-sum structure of (6), several stochastic optimization methods such as SAG [11], SDCA [12] and SVRG [9], provide linear convergence rates. All these algorithms can be readily applied to (6); we choose SVRG since it is memory efficient and easy to implement. We also

apply the recently developed accelerations techniques for first order optimization methods [13, 14] to obtain an accelerated SVRG (ASVRG) algorithm. We give the sketch of SVRG for (6) in Appendix C.

2 PN 2 Note that $f(u) = N \sum_{i=1}^N f_i(u)$ where each component $f_i(u) = \frac{1}{2} \|u - x_i\|_V^2 + \frac{\lambda}{2} \|u\|_X^2$ is k_i -smooth, and $f(u)$ is μ -strongly convex with $\mu = \min_i \mu_i$. We show in Appendix D that the initial suboptimality for minimizing $f(u)$ is upper-bounded by constant when using the warm-starts. We quote the convergence rates of SVRG [9] and ASVRG [14] below. ϵ satisfying $E[f(u) - f^*] \leq \epsilon$ Lemma 3. The SVRG algorithm [9] finds a vector u such that $E[f(u) - f^*] \leq \epsilon$ in time

$2 \max_i k_i \frac{1}{\mu} O_{\text{dx}}(N + \frac{1}{\mu}) \log \frac{1}{\epsilon}$ where $\frac{1}{\mu} = \frac{1}{\min_i \mu_i}$. The ASVRG algorithm [14] finds a such solution

in time $O_{\text{dx}}(N \frac{1}{\mu} \log \frac{1}{\epsilon})$. Remarks As mentioned in [14], the acceleration version provides speedup over normal SVRG only when $\frac{1}{\mu} \ll N$ and we only show the dominant term in the above complexity. By combining the iteration complexity of the outer loop (Theorem 2) and the time complexity of the inner loop (Lemma 3), we obtain the total time complexity of

$2 \frac{1}{\mu} \sum_{i=1}^N \frac{k_i}{\mu_i} \frac{1}{\mu} \log \frac{1}{\epsilon} \leq \frac{1}{\mu} (N + \frac{1}{\mu}) \log \frac{1}{\epsilon}$ for ALS+SVRG and $O(\frac{1}{\mu} \sum_{i=1}^N \frac{k_i}{\mu_i} \log \frac{1}{\epsilon})$ for ALS+ASVRG, where

$\frac{1}{\mu} := \max_i \frac{1}{\mu_i}$ and $O(\frac{1}{\mu}) = O(\frac{1}{\min_i \mu_i})$, $\frac{1}{\mu_i} = \frac{1}{\lambda + \mu_i}$ depends on λ and μ_i . Our algorithm does not require the initialization to be close to the optimum and converges globally. For comparison, the locally convergent AppGrad has a time complexity

$2 \frac{1}{\mu} \sum_{i=1}^N \frac{k_i}{\mu_i} \log \frac{1}{\epsilon}$, where $\frac{1}{\mu} := \max_i \frac{1}{\mu_i}$, $\frac{1}{\mu_i} = \frac{1}{\lambda + \mu_i}$. Note, [3, Theorem 2.1] of $O(\frac{1}{\mu} \sum_{i=1}^N \frac{k_i}{\mu_i} \log \frac{1}{\epsilon})$

in this complexity, the dataset size N and the least squares condition number $\frac{1}{\mu}$ are multiplied together because AppGrad essentially uses batch gradient descent as the least squares solver. Within our framework, we can use accelerated gradient descent (AGD, [15]) instead and obtain a globally

convergent algorithm with a total time complexity of $O(\frac{1}{\mu} \sum_{i=1}^N \frac{k_i}{\mu_i} \log \frac{1}{\epsilon})$.

3.2

2

Algorithm II: Shift-and-invert preconditioning (SI) with variance reduction

The second algorithm is inspired by the shift-and-invert preconditioning method for PCA [16, 17]. Instead of running power iterations on C as defined in (3), we will be running power iterations on $\frac{1}{\lambda} C$

$\frac{1}{\lambda} C = \frac{1}{\lambda} \sum_{i=1}^N \frac{1}{\mu_i} (x_i - x_i^*) (x_i - x_i^*)^T$, (7) $M = \frac{1}{\lambda} C = \frac{1}{\lambda} \sum_{i=1}^N \frac{1}{\mu_i} (x_i - x_i^*) (x_i - x_i^*)^T$

We omit the regularization in these constants, which are typically very small, to have concise expressions. The expectation is taken over random sampling of component functions. High probability error bounds can be obtained using the Markov's inequality. 6

5

1 and the number of calls to the least squares solver of $\text{ht}(\mathbf{u}, \mathbf{v})$ is $O(\log \frac{1}{\epsilon} \log \frac{1}{\delta} + \log \frac{1}{\delta})$.

3.2.2 Phase II: Final normalization In order to satisfy the CCA constraints, we perform a last normalization $\mathbf{q} = \frac{\mathbf{u}^T \mathbf{u}}{\|\mathbf{u}\|} \mathbf{u}, \mathbf{v} = \frac{\mathbf{v}^T \mathbf{v}}{\|\mathbf{v}\|} \mathbf{v}$.

(12)

\mathbf{q}, \mathbf{v} as our final approximate solution to (1). We show that this step does not cause And we output (\mathbf{u}, \mathbf{v}) much loss in the alignments, as stated below (see its proof in Appendix G). Theorem 5 (Convergence of Algorithm 3, Phase II). Let Phase I of Algorithm 3 outputs $(\mathbf{u}^T, \mathbf{v}^T)$ to (1) such that that satisfy (11). Then after (12), we obtain an approximate \mathbf{u}, \mathbf{v} solution (\mathbf{u}, \mathbf{v}) such that $\langle \mathbf{u}, \mathbf{v} \rangle = 1$, $\min(\langle \mathbf{u}, \mathbf{v} \rangle, \langle \mathbf{u}, \mathbf{v} \rangle) \geq \frac{1}{2}$, and $\|\mathbf{u}\|, \|\mathbf{v}\| \leq 2$. **3.2.3 Time complexity** We have shown in Theorem 4 that Phase I only approximately solves a small number of instances of (9). The normalization steps (10) require computing the projections of the training set which are reused for computing batch gradients of (9). The final normalization (12) is done only once and costs $O(dN)$. Therefore, the time complexity of our algorithm mainly comes from solving the least squares problems (9) using SGD methods in a blackbox fashion. And the time complexity for SGD methods depends on the condition number of (9). Denote $\kappa = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$.

It is clear that

$\lambda_{\max}(\mathbf{Q}) \leq \lambda_{\max}(\mathbf{Q}_{xx}) + \lambda_{\max}(\mathbf{Q}_{yy})$, $\lambda_{\min}(\mathbf{Q}) \geq \min(\lambda_{\min}(\mathbf{Q}_{xx}), \lambda_{\min}(\mathbf{Q}_{yy}))$.

We have shown in the proof of Theorem 4 that

$\lambda_{\max}(\mathbf{Q}) \leq \frac{9}{c_1}$

$\lambda_{\min}(\mathbf{Q}) \geq \frac{1}{c_1}$

$\lambda_{\max}(\mathbf{Q}_{xx}) \leq \frac{9}{c_1}$

$\lambda_{\min}(\mathbf{Q}_{xx}) \geq \frac{1}{c_1}$

$\lambda_{\max}(\mathbf{Q}_{yy}) \leq \frac{9}{c_1}$

Lemma 10, Appendix F.2), and thus the condition number for AGD is $\kappa = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})} \leq \frac{\lambda_{\max}(\mathbf{Q}_{xx}) + \lambda_{\max}(\mathbf{Q}_{yy})}{\min(\lambda_{\min}(\mathbf{Q}_{xx}), \lambda_{\min}(\mathbf{Q}_{yy}))}$.

throughout Algorithm 3 (cf. $\lambda_{\max}(\mathbf{Q}) \leq \frac{9}{c_1}$, $\lambda_{\min}(\mathbf{Q}) \geq \frac{1}{c_1}$).

$\lambda_{\max}(\mathbf{Q}_{xx}) \leq \frac{9}{c_1}$

$\lambda_{\min}(\mathbf{Q}_{xx}) \geq \frac{1}{c_1}$

where

$\kappa :=$ For SVRG/ASVRG, the relevant condition number depends on the gradient Lipschitz constant of individual components. We show in Appendix H (Lemma 12) that the maximum relevant condition number is at most $\frac{9}{c}$. An interesting issue for SVRG/ASVRG is that, depending on the value of ϵ , the independent components $\text{hit}(\mathbf{u}, \mathbf{v})$ may be nonconvex. If $\epsilon \leq 1$, each component is still guaranteed to be convex; otherwise, some PN components might be non-convex, with the overall average $\frac{1}{N} \sum_{i=1}^N \text{hit}_i$ being convex. In the later case, we use the modified analysis of SVRG [16, Appendix B] for its time

complexity. We use warmstart in SI as in ALS, and the initial suboptimality for each subproblem can be bounded similarly.

The total time complexities of our SI meta-algorithm are given in Table 1. Note that κ (or κ^2) are multiplied together, giving the effective condition number. When using SVRG as and κ^2

the least squares solver, we obtain the total $\kappa^2 \log \kappa + \kappa$ time complexity of $O(d(N + \kappa^2) \log \kappa)$ otherwise. When using all components are convex, and $O(\kappa^2)$

$\kappa^2 \log \kappa$ if all components are convex, and using ASVRG, we have $O(\kappa^2)$

hides poly-logarithmic dependences on dN otherwise. Here $O(\kappa^2 \log \kappa)$ is remarkable that the SI meta-algorithm is able to separate the dependence of dataset size and N from other parameters in the time complexities.

In a parallel work [6], the authors independently proposed a similar ALS algorithm,

and they solve the least squares problems using AGD. The time complexity of their algorithm for $dN \log \kappa$, which has linear dependence on the first canonical correlation is $O(\kappa^2)$

(so their algorithm is linearly convergent, but our complexity for ALS+AGD has on $\kappa^2 \log \kappa$)

quadratic dependence on this factor), but typically worse dependence on N and κ (see remarks in Section 3.1.1). Moreover, our SI algorithm tends to significantly outperform ALS theoretically and

empirically. It is future work to remove extra $\log \kappa$ dependence in our analysis.

Our arxiv preprint for the ALS meta-algorithm was posted before their paper got accepted by ICML 2016.

$\kappa_x = \kappa_y = 10^5$, $\kappa = 53340$, $\kappa = 5.345$
 $\kappa_x = \kappa_y = 10^4$, $\kappa = 5335$, $\kappa = 4.924$

0
Suboptimality
Median
10
CCALin SI-AVR
-5
10-5
10 SI-AVR
ALS-AVR
ALS-VR
10-4
AppGrad
-5
10 SI-AVR

SI-VR
 ALS-VR
 ALS-AVR
 10-10
 10-10
 -10
 10
 ALS-VR
 ALS-VR
 100
 SI-VR
 10-15
 100
 200
 300
 400
 500
 600
 0
 100
 200
 300
 SI-VR
 10-15
 400
 500
 0
 100
 200
 -15
 300
 400
 500
 600
 0
 $?? = 34070, ? = 10.58$
 100
 200
 300
 400
 500
 600
 $?? = 3416, ? = 9.082$ 0
 S-AppGrad 10
 100CCALin
 AppGrad

ALS-VR
 AppGrad
 AppGrad
 10-1
 ALS-AVR
 10
 0
 600
 $?? = 332800, ? = 11.10$ 10
 CCALin
 CCALin
 SI-VR
 ALS-AVR
 0
 JW11
 S-AppGrad
 S-AppGrad
 SI-AVR
 -2
 10
 -6
 Suboptimality
 AppGrad CCALin
 $?? = 2699000, ? = 11.22$
 CCALin
 S-AppGrad
 10-2
 S-AppGrad
 ALS-VR
 SI-AVR -2
 ALS-AVR
 10-5
 10
 S-AppGrad
 SI-VR
 10
 SI-AVR
 SI-AVR
 ALS-VR
 10-4 SI-VR
 -3
 ALS-AVR
 -5
 10
 CCALin
 AppGrad

ALS-AVR
 SI-AVR
 -10
 10
 ALS-VR ALS-AVR
 10-4
 10-6 0
 100
 200
 300
 400
 500
 600
 100
 SI-VR
 10-10
 0
 ?? = 2235000, ? = 12.82
 100
 200
 300
 400
 500
 600
 100
 100
 200
 300
 400
 500
 ?? = 22350, ? = 12.30
 100
 200
 300
 400
 500
 600
 ?? = 2236, ? = 9.874 0
 10
 ALS-VR
 AppGrad
 S-AppGrad
 AppGrad
 AppGrad
 S-AppGrad
 CCALin

S-AppGrad
 ALS-AVR
 0
 600
 100
 CCALin
 SI-VR -15
 10 0
 $?? = 223500, ? = 12.75$
 AppGrad
 Suboptimality
 0
 10 S-AppGrad
 AppGrad
 AppGrad
 10
 MNIST
 $?x = ?y = 10?2 ?? = 54.34, ? = 2.548$
 0 10 S-AppGrad
 0
 10CCALin
 $?x = ?y = 10?3 ?? = 534.4, ? = 4.256$
 CCALin ALS-VR
 10-2
 CCALin
 S-AppGrad
 ALS-AVR
 ALS-VR
 -5
 -5
 10
 -5
 10
 ALS-AVR
 10
 ALS-AVR
 SI-AVR
 10-4 ALS-VR
 -10
 10-10
 -10
 10
 10 SI-VR
 SI-VR
 10-6
 100

200
 300
 400
 # Passes
 500
 600
 10-15
 SI-AVR
 10-15 0
 100
 200
 300
 400
 500
 600
 # Passes
 0
 100
 SI-AVR
 SI-VR
 SI-VR
 SI-AVR
 0
 200
 300
 400
 # Passes
 500
 600
 10-15 0
 100
 200
 300
 400
 500
 600
 # Passes

Figure 1: Comparison of suboptimality vs. # passes for For each dataset and different algorithms. $\lambda_1(\lambda_{xx}) \lambda_{\max}(\lambda_{yy})$, regularization parameters (λ_x, λ_y) , we give $\lambda = \max \lambda_{\max} \lambda_{\min}(\lambda_{yy})$ and $\lambda = \lambda_2 \lambda_2 \cdot \min(\lambda_{xx}) 1$

Extension to multi-dimensional projections To extend our algorithms to L-dimensional projections, we can extract the dimensions sequentially and remove the explained correlation from λ_{xy} each time we extract a new dimension [18]. For the ALS meta-algorithm, a cleaner approach is to extract the L dimensions

simultaneously using (inexact) orthogonal iterations [8], in which case the sub-problems become multi-dimensional regressions and our normalization steps are of the form $\frac{1}{\sqrt{\lambda}} U^T \tilde{U} \tilde{U}^T U$ (the same normalization is used by [3, 4]). Such normalization involves $\frac{1}{\sqrt{\lambda}} (U^T \tilde{U} \tilde{U}^T U + U^T U)$ the eigenvalue decomposition of a $L \times L$ matrix and can be solved exactly as we typically look for low dimensional projections. Our analysis for $L = 1$ can be extended to this scenario and the convergence rate of ALS will depend on the gap between λ_L and λ_{L+1} .

4

Experiments

We demonstrate the proposed algorithms, namely ALS-VR, ALS-AVR, SI-VR, and SI-AVR, abbreviated as \tilde{U} meta-algorithm \tilde{U} least squares solver (VR for SVRG, and AVR for ASVRG) on three real-world datasets: Mediamill [19] ($N = 3 \times 10^4$), JW11 [20] ($N = 3 \times 10^4$), and MNIST [21] ($N = 6 \times 10^4$). We compare our algorithms with batch AppGrad and its stochastic version s-AppGrad [3], as well as the CCALin algorithm in parallel work [6]. For each algorithm, we compare the canonical correlation estimated by the iterates at different number of passes over the data with that of the exact solution by SVD. For each dataset, we vary the regularization parameters $\lambda_x = \lambda_y$ over $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$ to vary the least squares condition numbers, and larger regularization leads to better conditioning. We plot the suboptimality in objective vs. # passes for each algorithm in Figure 1. Experimental details (e.g. SVRG parameters) are given in Appendix I. We make the following observations from the results. First, the proposed stochastic algorithms significantly outperform batch gradient based methods AppGrad/CCALin. This is because the least squares condition numbers for these datasets are large, and SVRG enable us to decouple dependences on the dataset size N and the condition number κ in the time complexity. Second, SI-VR converges faster than ALS-VR as it further decouples the dependence on N and the singular value gap of T . Third, inexact normalizations keep the s-AppGrad algorithm from converging to an accurate solution. Finally, ASVRG improves over SVRG when the condition number is large. Acknowledgments Research partially supported by NSF BIGDATA grant 1546500. 8

2 References

- [1] H. Hotelling. Relations between two sets of variates. *Biometrika*, 28(3/4):321-377, 1936.
- [2] H. D. Vinod. Canonical ridge and econometrics of joint production. *J. Econometrics*, 1976.
- [3] Z. Ma, Y. Lu, and D. Foster. Finding linear structure in large datasets with scalable canonical correlation analysis. In *ICML*, 2015.
- [4] W. Wang, R. Arora, N. Srebro, and K. Livescu. Stochastic optimization for deep CCA via nonlinear orthogonal iterations. In *ALLERTON*, 2015.
- [5] B. Xie, Y. Liang, and L. Song. Scale up nonlinear component analysis with doubly stochastic gradients. In *NIPS*, 2015.
- [6] R. Ge, C. Jin, S. Kakade, P. Netrapalli, and A. Sidford. Efficient algorithms for large-scale generalized eigenvector computation and canonical correlation analysis. *arXiv*, April 13 2016.
- [7] G. Golub

and H. Zha. Linear Algebra for Signal Processing, chapter The Canonical Correlations of Matrix Pairs and their Numerical Computation, pages 27–49. 1995. [8] G. Golub and C. van Loan. Matrix Computations. third edition, 1996. [9] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In NIPS, 2013. [10] Y. Lu and D. Foster. Large scale canonical correlation analysis with iterative least squares. In NIPS, 2014. [11] M. Schmidt, N. Le Roux, and F. Bach. Minimizing ℓ_1 sums with the stochastic average gradient. Technical Report HAL 00860051, Ecole Normale Supérieure, 2013. [12] S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss minimization. Journal of Machine Learning Research, 2013. [13] R. Frostig, R. Ge, S. Kakade, and A. Sidford. Un-regularizing: Approximate proximal point and faster stochastic algorithms for empirical risk minimization. In ICML, 2015. [14] H. Lin, J. Mairal, and Z. Harchaoui. A universal catalyst for first-order optimization. In NIPS, 2015. [15] Y. Nesterov. Introductory Lectures on Convex Optimization. A Basic Course. Springer, 2004. [16] D. Garber and E. Hazan. Fast and simple PCA via convex optimization. arXiv, 2015. [17] C. Jin, S. Kakade, C. Musco, P. Netrapalli, and A. Sidford. Robust shift-and-invert preconditioning: Faster and more sample efficient algorithms for eigenvector computation. 2015. [18] D. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. Biostatistics, 2009. [19] C. Snoek, M. Worring, J. van Gemert, J. Geusebroek, and A. Smeulders. The challenge problem for automated detection of 101 semantic concepts in multimedia. In MULTIMEDIA, 2006. [20] J. Westbury. X-Ray Microbeam Speech Production Database User’s Handbook, 1994. [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. Proc. IEEE, 86(11):2278–2324, 1998. [22] M. Warmuth and D. Kuzmin. Randomized online PCA algorithms with regret bounds that are logarithmic in the dimension. Journal of Machine Learning Research, 2008. [23] R. Arora, A. Cotter, K. Livescu, and N. Srebro. Stochastic optimization for PCA and PLS. In ALLERTON, 2012. [24] A. Balsubramani, S. Dasgupta, and Y. Freund. The fast convergence of incremental PCA. In NIPS, 2013. [25] O. Shamir. A stochastic PCA and SVD algorithm with an exponential convergence rate. In ICML, 2015. [26] F. Yger, M. Berar, G. Gasso, and A. Rakotomamonjy. Adaptive canonical correlation analysis based on matrix manifolds. In ICML, 2012.