



Efficient Globally Convergent Stochastic Optimization for Canonical Correlation Analysis

Weiran Wang^{‡*}

Jialei Wang^{‡*}

Dan Garber[‡]

Nathan Srebro[‡]

[‡]Toyota Technological Institute at Chicago

[†]University of Chicago



Canonical correlation analysis (CCA)

Given data: $\mathbf{X} \in \mathbb{R}^{d_x \times N}$, $\mathbf{Y} \in \mathbb{R}^{d_y \times N}$, CCA solves

$$\begin{aligned} \max_{\mathbf{u}, \mathbf{v}} \quad & \mathbf{u}^\top \Sigma_{xy} \mathbf{v}, \\ \text{s.t.} \quad & \mathbf{u}^\top \Sigma_{xx} \mathbf{u} = \mathbf{v}^\top \Sigma_{yy} \mathbf{v} = 1, \end{aligned}$$

where $\Sigma_{xy} = \frac{1}{N} \mathbf{X} \mathbf{Y}^\top$, $\Sigma_{xx} = \frac{1}{N} \mathbf{X} \mathbf{X}^\top + \gamma_x \mathbf{I}$, $\Sigma_{yy} = \frac{1}{N} \mathbf{Y} \mathbf{Y}^\top + \gamma_y \mathbf{I}$, $(\gamma_x, \gamma_y) \geq 0$.

Goal: efficiently solve CCA problem when N, d are large.

Our approach: reduce CCA to a sequence of least square problems, and solve each problem to sufficient accuracy using fast solver.

Closed-form Solution

Form matrix \mathbf{T} :

$$\mathbf{T} = \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \in \mathbb{R}^{d_x \times d_y}$$

Compute SVD:

$$\mathbf{T} = [\boldsymbol{\phi}, \mathbf{a}_2, \dots, \mathbf{a}_r] \begin{bmatrix} \rho_1 & & & \\ & \rho_2 & & \\ & & \dots & \\ & & & \rho_r \end{bmatrix} [\boldsymbol{\psi}, \mathbf{b}_2, \dots, \mathbf{b}_r]^\top$$

Obtain $(\mathbf{u}^*, \mathbf{v}^*) = (\Sigma_{xx}^{-\frac{1}{2}} \boldsymbol{\phi}, \Sigma_{yy}^{-\frac{1}{2}} \boldsymbol{\psi})$

Forming \mathbf{T} costs $\mathcal{O}(Nd^2 + d^3)$

Alternating least squares (ALS) [Golub and Zha 1995]

The ALS Algorithm: For $t = 1, \dots$,

$$\begin{aligned} \tilde{\mathbf{u}}_t &\leftarrow \Sigma_{xx}^{-1} \Sigma_{xy} \mathbf{v}_{t-1}, & \tilde{\mathbf{v}}_t &\leftarrow \Sigma_{yy}^{-1} \Sigma_{xy}^\top \mathbf{u}_{t-1} \\ \mathbf{u}_t &\leftarrow \tilde{\mathbf{u}}_t / \sqrt{\tilde{\mathbf{u}}_t^\top \Sigma_{xx} \tilde{\mathbf{u}}_t}, & \mathbf{v}_t &\leftarrow \tilde{\mathbf{v}}_t / \sqrt{\tilde{\mathbf{v}}_t^\top \Sigma_{yy} \tilde{\mathbf{v}}_t} \end{aligned}$$

Alternative view in $\boldsymbol{\phi}_t = \Sigma_{xx}^{-\frac{1}{2}} \mathbf{u}_t$ and $\boldsymbol{\psi}_t = \Sigma_{yy}^{-\frac{1}{2}} \mathbf{v}_t$:

$$\begin{aligned} \tilde{\boldsymbol{\phi}}_t &\leftarrow \Sigma_{xx}^{-\frac{1}{2}} \Sigma_{xy} \Sigma_{yy}^{-\frac{1}{2}} \boldsymbol{\psi}_{t-1}, & \tilde{\boldsymbol{\psi}}_t &\leftarrow \Sigma_{yy}^{-\frac{1}{2}} \Sigma_{xy}^\top \Sigma_{xx}^{-\frac{1}{2}} \boldsymbol{\phi}_{t-1} \\ \boldsymbol{\phi}_t &\leftarrow \tilde{\boldsymbol{\phi}}_t / \|\tilde{\boldsymbol{\phi}}_t\|, & \boldsymbol{\psi}_t &\leftarrow \tilde{\boldsymbol{\psi}}_t / \|\tilde{\boldsymbol{\psi}}_t\| \end{aligned}$$

ALS w.r.t $(\mathbf{u}, \mathbf{v}) \iff$ Power iterations w.r.t $(\boldsymbol{\phi}, \boldsymbol{\psi})$

$$\begin{bmatrix} \tilde{\boldsymbol{\phi}}_t \\ \tilde{\boldsymbol{\psi}}_t \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{0} & \mathbf{T} \\ \mathbf{T}^\top & \mathbf{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{\phi}_{t-1} \\ \boldsymbol{\psi}_{t-1} \end{bmatrix}, \quad \begin{bmatrix} \boldsymbol{\phi}_t \\ \boldsymbol{\psi}_t \end{bmatrix} \leftarrow \begin{bmatrix} \tilde{\boldsymbol{\phi}}_t / \|\tilde{\boldsymbol{\phi}}_t\| \\ \tilde{\boldsymbol{\psi}}_t / \|\tilde{\boldsymbol{\psi}}_t\| \end{bmatrix}$$

Every two steps of ALS perform one step power iteration on $\mathbf{T} \mathbf{T}^\top$ for the $\boldsymbol{\phi}$ -sequence, and $\mathbf{T}^\top \mathbf{T}$ for the $\boldsymbol{\psi}$ -sequence.

Theorem: After $T = \mathcal{O}\left(\frac{\rho_1^2}{\rho_1^2 - \rho_2^2} \log\left(\frac{1}{\eta}\right)\right)$ iterations, we have for all $t \geq T$:

$$\min((\mathbf{u}_t^\top \Sigma_{xx} \mathbf{u}^*)^2, (\mathbf{v}_t^\top \Sigma_{yy} \mathbf{v}^*)^2) \geq 1 - \eta, \quad \mathbf{u}_t^\top \Sigma_{xy} \mathbf{v}_t \geq \rho_1(1 - 2\eta).$$

Inexact ALS: inexact power iteration on $\begin{bmatrix} \mathbf{0} & \mathbf{T} \\ \mathbf{T}^\top & \mathbf{0} \end{bmatrix}$

Non-zero eigenvalues: $\rho_1 \geq \rho_2 \geq \dots \geq \rho_r \geq -\rho_r \geq \dots \geq -\rho_1$

with eigenvectors

$$\begin{bmatrix} \boldsymbol{\phi} \\ \boldsymbol{\psi} \end{bmatrix}, \begin{bmatrix} \mathbf{a}_2 \\ \mathbf{b}_2 \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{a}_r \\ \mathbf{b}_r \end{bmatrix}, \begin{bmatrix} \mathbf{a}_r \\ -\mathbf{b}_r \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{a}_2 \\ -\mathbf{b}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\phi} \\ -\boldsymbol{\psi} \end{bmatrix}$$

Relative eigengap: $\frac{\rho_1^2}{\rho_1^2 - \rho_2^2}$

Algorithm: For $t = 1, \dots, T$

$$\text{-- Approximately solve } \min_{\mathbf{u}} f_t(\mathbf{u}) := \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\mathbf{u}^\top \mathbf{x}_i - \mathbf{v}_{t-1}^\top \mathbf{y}_i)^2 + \frac{\gamma_x}{2} \|\mathbf{u}\|^2 \quad \text{s.t.}$$

$$f_t(\tilde{\mathbf{u}}_t) \leq \min_{\mathbf{u}} f_t(\mathbf{u}) + \epsilon$$

$$\text{-- Approximately solve } \min_{\mathbf{v}} g_t(\mathbf{v}) := \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\mathbf{v}^\top \mathbf{y}_i - \mathbf{u}_{t-1}^\top \mathbf{x}_i)^2 + \frac{\gamma_y}{2} \|\mathbf{v}\|^2 \quad \text{s.t.}$$

$$g_t(\tilde{\mathbf{v}}_t) \leq \min_{\mathbf{v}} g_t(\mathbf{v}) + \epsilon$$

$$\text{-- Normalization: } \mathbf{u}_t \leftarrow \tilde{\mathbf{u}}_t / \sqrt{\tilde{\mathbf{u}}_t^\top \Sigma_{xx} \tilde{\mathbf{u}}_t}, \quad \mathbf{v}_t \leftarrow \tilde{\mathbf{v}}_t / \sqrt{\tilde{\mathbf{v}}_t^\top \Sigma_{yy} \tilde{\mathbf{v}}_t}$$

Subproblem conditioning: κ' for NAG, κ for SVRG

Theorem: After $T = \mathcal{O}\left(\frac{\rho_1^2}{\rho_1^2 - \rho_2^2} \log\left(\frac{2}{\eta}\right)\right)$ iterations with $\epsilon(T) = \frac{\eta^2 \rho_2^2}{128} \left(\frac{(2\rho_1/\rho_2) - 1}{(2\rho_1/\rho_2)^T - 1}\right)^2$, we have $\mathbf{u}_T^\top \Sigma_{xx} \mathbf{u}_T = \mathbf{v}_T^\top \Sigma_{yy} \mathbf{v}_T = 1$ and $\min((\mathbf{u}_T^\top \Sigma_{xx} \mathbf{u}^*)^2, (\mathbf{v}_T^\top \Sigma_{yy} \mathbf{v}^*)^2) \geq 1 - \eta$.

Shift-and-invert: inexact power iteration on $\begin{bmatrix} \lambda \mathbf{I} & -\mathbf{T} \\ -\mathbf{T}^\top & \lambda \mathbf{I} \end{bmatrix}^{-1}$

Eigenvalues: $\frac{1}{\lambda - \rho_1} \geq \dots \geq \frac{1}{\lambda - \rho_r} \geq \dots \geq \frac{1}{\lambda + \rho_r} \geq \dots \geq \frac{1}{\lambda + \rho_1}$

with eigenvectors $\frac{1}{\sqrt{2}} \begin{bmatrix} \boldsymbol{\phi} \\ \boldsymbol{\psi} \end{bmatrix}, \dots, \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{a}_r \\ \mathbf{b}_r \end{bmatrix}, \dots, \frac{1}{\sqrt{2}} \begin{bmatrix} \mathbf{a}_r \\ -\mathbf{b}_r \end{bmatrix}, \dots, \frac{1}{\sqrt{2}} \begin{bmatrix} \boldsymbol{\phi} \\ -\boldsymbol{\psi} \end{bmatrix}$

For $\lambda = \rho_1 + c(\rho_1 - \rho_2)$ where $c = \mathcal{O}(1)$, relative eigengap is $1 + c$, a constant

Algorithm: For $t = 1, \dots, T$

$$\text{-- Approximately solve } \begin{bmatrix} \tilde{\boldsymbol{\phi}}_t \\ \tilde{\boldsymbol{\psi}}_t \end{bmatrix} \leftarrow \begin{bmatrix} \lambda \mathbf{I} & -\mathbf{T} \\ -\mathbf{T}^\top & \lambda \mathbf{I} \end{bmatrix}^{-1} \begin{bmatrix} \boldsymbol{\phi}_{t-1} \\ \boldsymbol{\psi}_{t-1} \end{bmatrix} \quad \text{or}$$

$$(\tilde{\mathbf{u}}_t, \tilde{\mathbf{v}}_t) \approx \arg \min_{\mathbf{u}, \mathbf{v}} \frac{1}{2} [\mathbf{u}^\top \mathbf{v}^\top] \begin{bmatrix} \lambda \Sigma_{xx} & -\Sigma_{xy} \\ -\Sigma_{xy}^\top & \lambda \Sigma_{yy} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} - \mathbf{u}^\top \Sigma_{xx} \mathbf{u}_{t-1} - \mathbf{v}^\top \Sigma_{yy} \mathbf{v}_{t-1}$$

$$\text{-- Intermediate normalization: } \begin{bmatrix} \tilde{\mathbf{u}}_t \\ \tilde{\mathbf{v}}_t \end{bmatrix} \leftarrow \sqrt{2} \begin{bmatrix} \tilde{\mathbf{u}}_t \\ \tilde{\mathbf{v}}_t \end{bmatrix} / \sqrt{\tilde{\mathbf{u}}_t^\top \Sigma_{xx} \tilde{\mathbf{u}}_t + \tilde{\mathbf{v}}_t^\top \Sigma_{yy} \tilde{\mathbf{v}}_t}$$

$$\text{Final normalization: } \hat{\mathbf{u}} \leftarrow \mathbf{u}_T / \sqrt{\mathbf{u}_T^\top \Sigma_{xx} \mathbf{u}_T}, \quad \hat{\mathbf{v}} \leftarrow \mathbf{v}_T / \sqrt{\mathbf{v}_T^\top \Sigma_{yy} \mathbf{v}_T}$$

Subproblem conditioning: $\mathcal{O}\left(\frac{\kappa'}{\rho_1 - \rho_2}\right)$ for NAG, $\mathcal{O}\left(\frac{\kappa}{\rho_1 - \rho_2}\right)$ for SVRG

Theorem: After $T = \mathcal{O}\left(\log\left(\frac{2}{\eta}\right)\right)$ iterations, we have $\hat{\mathbf{u}}^\top \Sigma_{xx} \hat{\mathbf{u}} = \hat{\mathbf{v}}^\top \Sigma_{yy} \hat{\mathbf{v}} = 1$, and $\min((\hat{\mathbf{u}}^\top \Sigma_{xx} \mathbf{u}^*)^2, (\hat{\mathbf{v}}^\top \Sigma_{yy} \mathbf{v}^*)^2) \geq 1 - \eta$.

There exists an efficient algorithm for locating λ .

Summary of running time

	Alternating least squares	Shift-and-invert ($\lambda > \rho_1$)
Power iterations	$\begin{bmatrix} \mathbf{0} & \mathbf{T} \\ \mathbf{T}^\top & \mathbf{0} \end{bmatrix}$	$\begin{bmatrix} \lambda \mathbf{I} & -\mathbf{T} \\ -\mathbf{T}^\top & \lambda \mathbf{I} \end{bmatrix}^{-1}$
Least squares	$\begin{bmatrix} \Sigma_{xx} & \\ & \Sigma_{yy} \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{xy} \mathbf{v}_{t-1} \\ \Sigma_{xy}^\top \mathbf{u}_{t-1} \end{bmatrix}$	$\begin{bmatrix} \lambda \Sigma_{xx} & -\Sigma_{xy} \\ -\Sigma_{xy}^\top & \lambda \Sigma_{yy} \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{xx} \mathbf{u}_{t-1} \\ \Sigma_{yy} \mathbf{v}_{t-1} \end{bmatrix}$
NAG	$dN \sqrt{\kappa'} \left(\frac{\rho_1^2}{\rho_1^2 - \rho_2^2}\right)^2 \cdot \log^2\left(\frac{1}{\eta}\right)$	$dN \sqrt{\tilde{\kappa}'} \sqrt{\frac{1}{\rho_1 - \rho_2}} \cdot \log^2\left(\frac{1}{\eta}\right)$
SVRG	$d(N + \kappa) \left(\frac{\rho_1^2}{\rho_1^2 - \rho_2^2}\right)^2 \cdot \log^2\left(\frac{1}{\eta}\right)$	$d \left(N + \left(\tilde{\kappa} \frac{1}{\rho_1 - \rho_2}\right)^2\right) \cdot \log^2\left(\frac{1}{\eta}\right)$
ASVRG	$d \sqrt{N} \sqrt{\kappa} \left(\frac{\rho_1^2}{\rho_1^2 - \rho_2^2}\right)^2 \cdot \log^2\left(\frac{1}{\eta}\right)$	$d N^{\frac{3}{2}} \sqrt{\tilde{\kappa}} \sqrt{\frac{1}{\rho_1 - \rho_2}} \cdot \log^2\left(\frac{1}{\eta}\right)$

Runtime (up to polylog factors) to find (\mathbf{u}, \mathbf{v}) with $\min((\mathbf{u}^\top \Sigma_{xx} \mathbf{u}^*)^2, (\mathbf{v}^\top \Sigma_{yy} \mathbf{v}^*)^2) \geq 1 - \eta$

$$\begin{aligned} \kappa' &:= \max\left(\frac{\sigma_{\max}(\Sigma_{xx}), \sigma_{\max}(\Sigma_{yy})}{\sigma_{\min}(\Sigma_{xx}), \sigma_{\min}(\Sigma_{yy})}\right) \leq \kappa := \max\left(\frac{\max_i \|\mathbf{x}_i\|^2, \max_i \|\mathbf{y}_i\|^2}{\sigma_{\min}(\Sigma_{xx}), \sigma_{\min}(\Sigma_{yy})}\right) \\ \wedge & \\ \tilde{\kappa}' &:= \frac{\max(\sigma_{\max}(\Sigma_{xx}), \sigma_{\max}(\Sigma_{yy}))}{\min(\sigma_{\min}(\Sigma_{xx}), \sigma_{\min}(\Sigma_{yy}))} \leq \tilde{\kappa} := \frac{\max_i \max(\|\mathbf{x}_i\|^2, \|\mathbf{y}_i\|^2)}{\min(\sigma_{\min}(\Sigma_{xx}), \sigma_{\min}(\Sigma_{yy}))} \end{aligned}$$

Related Work

- [Ma, Lu, and Foster 2015], AppGrad: one full gradient step on ALS subproblems, can only guarantee local convergence
- Parallel work: [Ge et al, 2016], CCAIn: ALS with subproblems solved by NAG, time complexity $\mathcal{O}\left(dN \sqrt{\kappa'} \frac{\rho_1^2}{\rho_1^2 - \rho_2^2} \cdot \log\left(\frac{1}{\eta}\right)\right)$

Experiments

Datasets	Description	d_x	d_y	N
JW11	Acoustic and articulation measurements	273	112	30,000
MNIST	Left and right halves of images	392	392	60,000

