

Towards Practical Alternating Least-Squares for CCA

Zhiqiang Xu Ping Li

Cognitive Computing Lab

Baidu Research

{xuzhiqiang04, liping11}@baidu.com

November 3, 2019

Canonical Correlation Analysis (CCA)

- $\mathbf{X} \in \mathbb{R}^{d_x \times n}$ and $\mathbf{Y} \in \mathbb{R}^{d_y \times n}$: given data matrix pair
- Empirical cross-covariance matrix

$$\mathbf{C}_{xy} = \frac{1}{n} \mathbf{X} \mathbf{Y}^\top$$

- Two empirical auto-covariance matrices

$$\mathbf{C}_{xx} = \frac{1}{n} \mathbf{X} \mathbf{X}^\top + r_x \mathbf{I}, \quad \mathbf{C}_{yy} = \frac{1}{n} \mathbf{Y} \mathbf{Y}^\top + r_y \mathbf{I}$$

CCA

$$\max_{\substack{\Phi^\top \mathbf{C}_{xx} \Phi = \mathbf{I} \\ \Psi^\top \mathbf{C}_{yy} \Psi = \mathbf{I}}} \text{tr}(\Phi^\top \mathbf{C}_{xy} \Psi)$$

where $\Phi \in \mathbb{R}^{d_x \times k}$ and $\Psi \in \mathbb{R}^{d_y \times k}$, $k \geq 1$

Canonical Correlation Analysis (CCA)

Optimal Solution

$$(\Phi^*, \Psi^*) = (\mathbf{C}_{xx}^{-\frac{1}{2}} \mathbf{P}, \mathbf{C}_{yy}^{-\frac{1}{2}} \mathbf{Q}) \triangleq (\mathbf{U}, \mathbf{V}),$$

- $\mathbf{P} \in \mathbb{R}^{d_x \times k}$ and $\mathbf{Q} \in \mathbb{R}^{d_y \times k}$: top- k left and right singular subspaces of $\mathbf{C} = \mathbf{C}_{xx}^{-\frac{1}{2}} \mathbf{C}_{xy} \mathbf{C}_{yy}^{-\frac{1}{2}}$ in Euclidean metrics, respectively,

$$\mathbf{C} \stackrel{\text{SVD}}{=} \mathbf{P} \mathbf{\Sigma} \mathbf{Q}^\top + \mathbf{P}_\perp \mathbf{\Sigma}_\perp \mathbf{Q}_\perp^\top,$$

where $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_k)$ with σ_i being the i -th largest singular value, and $\mathbf{\Sigma}_\perp \in \mathbb{R}^{(r-k) \times (r-k)}$ is diagonal with $r = \text{rank}(\mathbf{C})$.

- $\mathbf{U} \in \mathbb{R}^{d_x \times k}$ and $\mathbf{V} \in \mathbb{R}^{d_y \times k}$: top- k canonical subspaces, which are in metrics \mathbf{C}_{xx} and \mathbf{C}_{yy} , respectively,

$$\mathbf{C}_{xy} = \mathbf{C}_{xx} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{C}_{yy} + \mathbf{C}_{xx} \mathbf{U}_\perp \mathbf{\Sigma}_\perp \mathbf{V}_\perp^\top \mathbf{C}_{yy}$$

Alternating Least-Squares (ALS)

Generalized Eigenvalue Problem Formulation

$$\mathbf{A}\mathbf{w} = \lambda\mathbf{B}\mathbf{w} \quad (\text{cca-stationarity})$$

where

$$\mathbf{A} = \begin{pmatrix} & \mathbf{C}_{xy} \\ \mathbf{C}_{xy}^\top & \end{pmatrix} \quad \text{and} \quad \mathbf{B} = \begin{pmatrix} \mathbf{C}_{xx} & \\ & \mathbf{C}_{yy} \end{pmatrix}.$$

- (\mathbf{A}, \mathbf{B}) 's generalized eigenvectors

$$\frac{1}{\sqrt{2}} \left\{ \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{v}_1 \end{pmatrix}, \begin{pmatrix} \mathbf{u}_2 \\ \mathbf{v}_2 \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{u}_r \\ \mathbf{v}_r \end{pmatrix}, \begin{pmatrix} -\mathbf{u}_r \\ \mathbf{v}_r \end{pmatrix}, \begin{pmatrix} -\mathbf{u}_2 \\ \mathbf{v}_2 \end{pmatrix}, \dots, \begin{pmatrix} -\mathbf{u}_1 \\ \mathbf{v}_1 \end{pmatrix} \right\},$$

corresponding to generalized eigenvalues $\sigma_1, \sigma_2, \dots, \sigma_r, -\sigma_r, -\sigma_2, \dots, -\sigma_1$, respectively, where

$$\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k), \quad \mathbf{V} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k)$$

Alternating Least-Squares (ALS)

Power Method for GEV

$$\tilde{\Omega}_{t+1} = \mathbf{B}^{-1} \mathbf{A} \Omega_t + \varsigma_t, \quad \Omega_{t+1} = \tilde{\Omega}_{t+1} (\tilde{\Omega}_{t+1}^\top \mathbf{B} \tilde{\Omega}_{t+1})^{-\frac{1}{2}}$$

1

- Suppose $\Omega_t \in \mathbb{R}^{(d_x+d_y) \times j}$, where j is even.
- Column space of Ω_t converges to a top- j generalized eigenspace spanned by $\frac{1}{2} \left\{ \begin{pmatrix} \mathbf{u}_1 \\ \mathbf{v}_1 \end{pmatrix}, \begin{pmatrix} -\mathbf{u}_1 \\ \mathbf{v}_1 \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{u}_{\frac{j}{2}} \\ \mathbf{v}_{\frac{j}{2}} \end{pmatrix}, \begin{pmatrix} -\mathbf{u}_{\frac{j}{2}} \\ \mathbf{v}_{\frac{j}{2}} \end{pmatrix} \right\}$ corresponding to j largest generalized eigenvalues $\sigma_1, -\sigma_1, \dots, \sigma_{\frac{j}{2}}, -\sigma_{\frac{j}{2}}$ in magnitude.

¹Rong Ge et al. "Efficient Algorithms for Large-scale Generalized Eigenvector Computation and Canonical Correlation Analysis". In: *International Conference on Machine Learning*. 2016, pp. 2741–2750.

Alternating Least-Squares (ALS)

Update Equations

$$\begin{cases} \tilde{\Phi}_t = \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_{t-1} + \xi_{t-1}, & \Phi_t = \tilde{\Phi}_t (\tilde{\Phi}_t^\top \mathbf{C}_{xx} \tilde{\Phi}_t + \tilde{\Psi}_t^\top \mathbf{C}_{yy} \tilde{\Psi}_t)^{-\frac{1}{2}} \\ \tilde{\Psi}_t = \mathbf{C}_{yy}^{-1} \mathbf{C}_{yx}^\top \Phi_{t-1} + \eta_{t-1}, & \Psi_t = \tilde{\Psi}_t (\tilde{\Psi}_t^\top \mathbf{C}_{yy} \tilde{\Psi}_t + \tilde{\Phi}_t^\top \mathbf{C}_{xx} \tilde{\Phi}_t)^{-\frac{1}{2}} \end{cases}$$

- To recover \mathbf{U} and \mathbf{V} , one needs to first set $j = 2k$ for having $\text{col}(\Omega_\infty)$ spanned by $\frac{1}{2} \begin{pmatrix} \mathbf{U} & -\mathbf{U} \\ \mathbf{V} & \mathbf{V} \end{pmatrix}$
- that is, $\Phi_t \in \mathbb{R}^{d_x \times 2k}$ and $\Psi_t \in \mathbb{R}^{d_y \times 2k}$
- then do random projection with random Gaussian matrix $\mathbf{G} \in \mathbb{R}^{2k \times k}$,

$$\begin{cases} \tilde{\Phi}_T = \Phi_T \mathbf{G} \\ \tilde{\Psi}_T = \Psi_T \mathbf{G} \end{cases}, \quad \begin{cases} \hat{\mathbf{U}} = \tilde{\Phi}_T (\tilde{\Phi}_T^\top \mathbf{C}_{xx} \tilde{\Phi}_T)^{-\frac{1}{2}} \\ \hat{\mathbf{V}} = \tilde{\Psi}_T (\tilde{\Psi}_T^\top \mathbf{C}_{yy} \tilde{\Psi}_T)^{-\frac{1}{2}} \end{cases}$$

Truly Alternating Least-Squares

Coupled Equations of Half the Size

$$\begin{cases} \tilde{\Phi}_t = \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_{t-1} + \xi_{t-1}, & \Phi_t = \tilde{\Phi}_t (\tilde{\Phi}_t^\top \mathbf{C}_{xx} \tilde{\Phi}_t)^{-\frac{1}{2}} \\ \tilde{\Psi}_t = \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top \Phi_t + \eta_t, & \Psi_t = \tilde{\Psi}_t (\tilde{\Psi}_t^\top \mathbf{C}_{yy} \tilde{\Psi}_t)^{-\frac{1}{2}} \end{cases}$$

- $\Phi_t \in \mathbb{R}^{d_x \times k}$ and $\Psi_t \in \mathbb{R}^{d_y \times k}$, more memory-efficient
- No need of the random projection any more

$$\hat{\mathbf{U}} = \Phi_T \text{ and } \hat{\mathbf{V}} = \Psi_T$$

- Faster roughly by a factor of $\frac{\sigma_k}{\sigma_k + \sigma_{k+1}}$, especially for cases of a small singular value gap

Truly Alternating Least-Squares (TALS)

TALS-CCA

Algorithm 1: TALS-CCA

- 1: **Input:** T, k , data matrices \mathbf{X}, \mathbf{Y}
 - 2: **Output:** approximate top- k canonical subspaces (Φ_T, Ψ_T)
 - 3: $\Phi_0 = \text{GS}_{\mathbf{C}_{xx}}(\Phi_{\text{init}})$, $\Phi_{\text{init}} \in \mathbb{R}^{d_x \times k}$ is random Gaussian
 $\Psi_0 = \text{GS}_{\mathbf{C}_{yy}}(\Psi_{\text{init}})$, $\Psi_{\text{init}} \in \mathbb{R}^{d_y \times k}$ is random Gaussian
 - 4: **for** $t = 1, 2, \dots, T$ **do**
 - 5: $\tilde{\Phi}_t \approx \arg \min \ell_t(\tilde{\Phi}) = \frac{1}{2n} \|\mathbf{X}^\top \tilde{\Phi} - \mathbf{Y}^\top \Psi_{t-1}\|_F^2 + \frac{r_x}{2} \|\tilde{\Phi}\|_F^2$
 starting from $\tilde{\Phi}^{(0)} = \Phi_{t-1}(\Phi_{t-1}^\top \mathbf{C}_{xx} \Phi_{t-1})^{-1}(\Phi_{t-1}^\top \mathbf{C}_{xy} \Psi_{t-1})$
 - 6: $\Phi_t = \text{GS}_{\mathbf{C}_{xx}}(\tilde{\Phi}_t)$
 - 7: $\tilde{\Psi}_t \approx \arg \min s_t(\tilde{\Psi}) = \frac{1}{2n} \|\mathbf{Y}^\top \tilde{\Psi} - \mathbf{X}^\top \Phi_t\|_F^2 + \frac{r_y}{2} \|\tilde{\Psi}\|_F^2$
 starting from $\tilde{\Psi}^{(0)} = \Psi_{t-1}(\Psi_{t-1}^\top \mathbf{C}_{yy} \Psi_{t-1})^{-1}(\Psi_{t-1}^\top \mathbf{C}_{xy}^\top \Phi_t)$
 - 8: $\Psi_t = \text{GS}_{\mathbf{C}_{yy}}(\tilde{\Psi}_t)$
 - 9: **end for**
-

Faster Alternating Least-Squares (FALS)

Momentum Acceleration

$$\begin{cases} \tilde{\Phi}_t = \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_{t-1} - \beta \Phi_{t-2} + \xi_{t-1}, & \Phi_t = \tilde{\Phi}_t (\tilde{\Phi}_t^\top \mathbf{C}_{xx} \tilde{\Phi}_t)^{-\frac{1}{2}} \\ \tilde{\Psi}_t = \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^\top \Phi_t - \beta \Psi_{t-1} + \eta_t, & \Psi_t = \tilde{\Psi}_t (\tilde{\Psi}_t^\top \mathbf{C}_{yy} \tilde{\Psi}_t)^{-\frac{1}{2}} \end{cases}$$

- Memory-efficient: $\Phi_t \in \mathbb{R}^{d_x \times k}$ and $\Psi_t \in \mathbb{R}^{d_y \times k}$
- No need of the random projection: $\hat{\mathbf{U}} = \Phi_T$ and $\hat{\mathbf{V}} = \Psi_T$
- Least-squares becomes

$$\begin{cases} \min l_t(\tilde{\Phi}) = \frac{1}{2n} \|\mathbf{X}^\top (\tilde{\Phi} + \beta \Phi_{t-2}) - \mathbf{Y}^\top \Psi_{t-1}\|_F^2 + \frac{r_x}{2} \|\tilde{\Phi} + \beta \Phi_{t-2}\|_F^2, \\ \min s_t(\tilde{\Psi}) = \frac{1}{2n} \|\mathbf{Y}^\top (\tilde{\Psi} + \beta \Psi_{t-1}) - \mathbf{X}^\top \Phi_t\|_F^2 + \frac{r_y}{2} \|\tilde{\Psi} + \beta \Psi_{t-1}\|_F^2 \end{cases}$$

Faster Alternating Least-Squares (FALS)

FALS-CCA

Algorithm 2: FALS-CCA

- 1: **Input:** T, k , momentum parameter β , data matrices \mathbf{X}, \mathbf{Y}
 - 2: **Output:** approximate top- k canonical subspaces (Φ_T, Ψ_T)
 - 3: $\Phi_{-1} = \mathbf{0} \in \mathbb{R}^{d_x \times k}$
 $\Phi_0 = \text{GS}_{\mathbf{C}_{xx}}(\Phi_{\text{init}})$, $\Phi_{\text{init}} \in \mathbb{R}^{d_x \times k}$ is random Gaussian
 $\Psi_0 = \text{GS}_{\mathbf{C}_{yy}}(\Psi_{\text{init}})$, $\Psi_{\text{init}} \in \mathbb{R}^{d_y \times k}$ is random Gaussian
 - 4: **for** $t = 1, 2, \dots, T$ **do**
 - 5: $\tilde{\Phi}_t \approx \arg \min l_t(\tilde{\Phi}) = \frac{1}{2n} \|\mathbf{X}^\top (\tilde{\Phi} + \beta \Phi_{t-2}) - \mathbf{Y}^\top \Psi_{t-1}\|_F^2 + \frac{r_x}{2} \|\tilde{\Phi} + \beta \Phi_{t-2}\|_F^2$
 starting from $\tilde{\Phi}^{(0)} = \Phi_{t-1} (\Phi_{t-1}^\top \mathbf{C}_{xx} \Phi_{t-1})^{-1} (\Phi_{t-1}^\top \mathbf{C}_{xy} \Psi_{t-1})$
 - 6: $\Phi_t = \text{GS}_{\mathbf{C}_{xx}}(\tilde{\Phi}_t)$
 - 7: $\tilde{\Psi}_t \approx \arg \min s_t(\tilde{\Psi}) = \frac{1}{2n} \|\mathbf{Y}^\top (\tilde{\Psi} + \beta \Psi_{t-1}) - \mathbf{X}^\top \Phi_t\|_F^2 + \frac{r_y}{2} \|\tilde{\Psi} + \beta \Psi_{t-1}\|_F^2$
 starting from $\tilde{\Psi}^{(0)} = \Psi_{t-1} (\Psi_{t-1}^\top \mathbf{C}_{yy} \Psi_{t-1})^{-1} (\Psi_{t-1}^\top \mathbf{C}_{xy}^\top \Phi_t)$
 - 8: $\Psi_t = \text{GS}_{\mathbf{C}_{yy}}(\tilde{\Psi}_t)$
 - 9: **end for**
-

Faster Alternating Least-Squares (FALS- T_{tals})

- Optimal² momentum parameter β should be around $\frac{\sigma_{k+1}^2}{4}$
- Σ estimates by $(\mathbf{U}^\top \mathbf{C}_{xx} \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{C}_{xy} \mathbf{V} = \Sigma = (\mathbf{V}^\top \mathbf{C}_{yy} \mathbf{V})^{-1} \mathbf{V}^\top \mathbf{C}_{xy}^\top \mathbf{U}$:

$$\Sigma^{(t,1)} = (\Phi_t^\top \mathbf{C}_{xx} \Phi_t)^{-1} \Phi_t^\top \mathbf{C}_{xy} \Psi_t,$$

$$\Sigma^{(t,2)} = (\Psi_{t-1}^\top \mathbf{C}_{yy} \Psi_{t-1})^{-1} \Psi_{t-1}^\top \mathbf{C}_{xy}^\top \Phi_t$$

- β estimates by TALS:

$$\hat{\beta} = \frac{1}{4} \min_i (\Sigma_{ii}^{(T_{tals},1)})^2 \text{ or } \frac{1}{4} \min_i (\Sigma_{ii}^{(T_{tals},2)})^2,$$

where T_{tals} is small.

²Peng Xu et al. "Accelerated Stochastic Power Iteration". In: *International Conference on Artificial Intelligence and Statistics, AISTATS 2018, 9-11 April 2018, Playa Blanca, Lanzarote, Canary Islands, Spain*. 2018, pp. 58–67. URL: <http://proceedings.mlr.press/v84/xu18a.html>.

Adaptive Alternating Least-Squares (AALS)

Adaptive Momentum Acceleration

$$\begin{cases} \tilde{\Phi}_t = \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \Psi_{t-1} - \beta_{t-1}^{\phi} \Phi_{t-2} + \xi_{t-1}, & \Phi_t = \tilde{\Phi}_t (\tilde{\Phi}_t^{\top} \mathbf{C}_{xx} \tilde{\Phi}_t)^{-\frac{1}{2}} \\ \tilde{\Psi}_t = \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^{\top} \Phi_t - \beta_t^{\psi} \Psi_{t-1} + \eta_t, & \Psi_t = \tilde{\Psi}_t (\tilde{\Psi}_t^{\top} \mathbf{C}_{yy} \tilde{\Psi}_t)^{-\frac{1}{2}} \end{cases}$$

- Adaptive momentum parameters:

$$\beta_t^{\phi} = \frac{1}{4} \min_i (\Sigma_{ii}^{(t,1)})^2 \text{ and } \beta_t^{\psi} = \frac{1}{4} \min_i (\Sigma_{ii}^{(t,2)})^2$$

- Least-squares becomes

$$\begin{cases} \min l_t(\tilde{\Phi}) = \frac{1}{2n} \|\mathbf{X}^{\top} (\tilde{\Phi} + \beta_{t-1}^{\phi} \Phi_{t-2}) - \mathbf{Y}^{\top} \Psi_{t-1}\|_F^2 + \frac{r_x}{2} \|\tilde{\Phi} + \beta_{t-1}^{\phi} \Phi_{t-2}\|_F^2, \\ \min s_t(\tilde{\Psi}) = \frac{1}{2n} \|\mathbf{Y}^{\top} (\tilde{\Psi} + \beta_t^{\psi} \Psi_{t-1}) - \mathbf{X}^{\top} \Phi_t\|_F^2 + \frac{r_y}{2} \|\tilde{\Psi} + \beta_t^{\psi} \Psi_{t-1}\|_F^2, \end{cases}$$

with initials $\tilde{\Phi}^{(0)} = \Phi_{t-1} \Sigma^{(t-1,1)}$, $\tilde{\Psi}^{(0)} = \Psi_{t-1} \Sigma^{(t-1,2)}$, respectively

Adaptive Alternating Least-Squares (AALS)

AALS-CCA

Algorithm 3: AALS-CCA

- 1: **Input:** T, k , data matrices \mathbf{X}, \mathbf{Y}
 - 2: **Output:** approximate top- k canonical subspaces (Φ_T, Ψ_T)
 - 3: $\Phi_{-1} = \mathbf{0} \in \mathbb{R}^{d_x \times k}$
 $\Phi_0 = \text{GSC}_{xx}(\Phi_{\text{init}})$, $\Phi_{\text{init}} \in \mathbb{R}^{d_x \times k}$ is random Gaussian
 $\Psi_0 = \text{GSC}_{yy}(\Psi_{\text{init}})$, $\Psi_{\text{init}} \in \mathbb{R}^{d_y \times k}$ is random Gaussian
 - 4: **for** $t = 1, 2, \dots, T$ **do**
 - 5: $\beta_{t-1}^\phi = \frac{1}{4} \min_i (\Sigma_{ii}^{(t-1,1)})^2$, where $\Sigma^{(t-1,1)} = (\Phi_{t-1}^\top \mathbf{C}_{xx} \Phi_{t-1})^{-1} \Phi_{t-1}^\top \mathbf{C}_{xy} \Psi_{t-1}$
 - 6: $\tilde{\Phi}_t \approx \arg \min l_t(\tilde{\Phi}) = \frac{1}{2n} \|\mathbf{X}^\top (\tilde{\Phi} + \beta_{t-1}^\phi \Phi_{t-2}) - \mathbf{Y}^\top \Psi_{t-1}\|_F^2 + \frac{r_x}{2} \|\tilde{\Phi} + \beta_{t-1}^\phi \Phi_{t-2}\|_F^2$
starting from $\tilde{\Phi}^{(0)} = \Phi_{t-1} \Sigma^{(t-1,1)}$
 - 7: $\Phi_t = \text{GSC}_{xx}(\tilde{\Phi}_t)$
 - 8: $\beta_t^\psi = \frac{1}{4} \min_i (\Sigma_{ii}^{(t,2)})^2$, where $\Sigma^{(t,2)} = (\Psi_{t-1}^\top \mathbf{C}_{yy} \Psi_{t-1})^{-1} \Psi_{t-1}^\top \mathbf{C}_{xy} \Phi_t$
 - 9: $\tilde{\Psi}_t \approx \arg \min s_t(\tilde{\Psi}) = \frac{1}{2n} \|\mathbf{Y}^\top (\tilde{\Psi} + \beta_t^\psi \Psi_{t-1}) - \mathbf{X}^\top \Phi_t\|_F^2 + \frac{r_y}{2} \|\tilde{\Psi} + \beta_t^\psi \Psi_{t-1}\|_F^2$
starting from $\tilde{\Psi}^{(0)} = \Psi_{t-1} \Sigma^{(t,2)}$
 - 10: $\Psi_t = \text{GSC}_{yy}(\tilde{\Psi}_t)$
 - 11: **end for**
-

Experiments

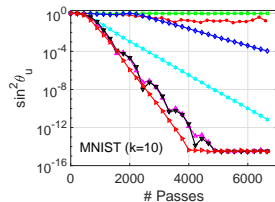
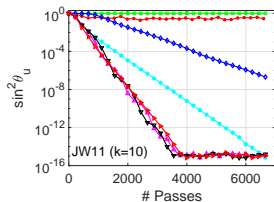
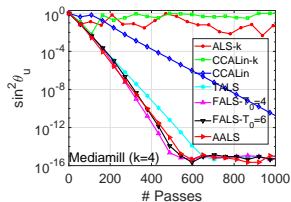
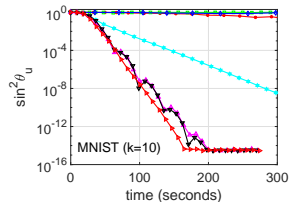
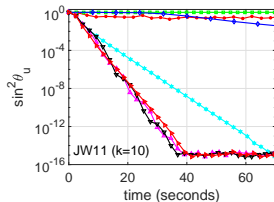
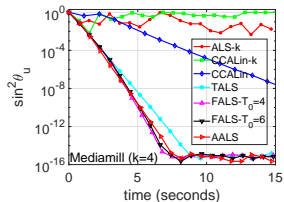
- Data

DATA	Description	d_x	d_y	n
Memdiamill	images and its labels	100	120	30000
JW11	acoustic and articulation	273	112	30000
MNIST	left and right halves of images	392	392	60000

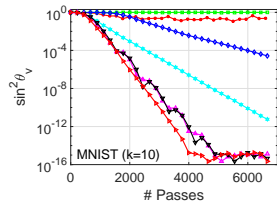
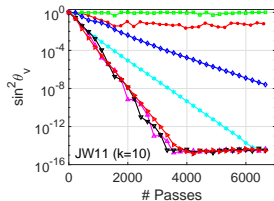
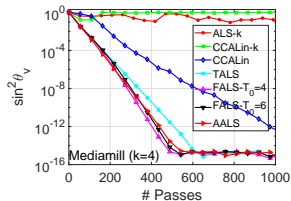
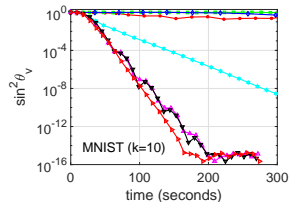
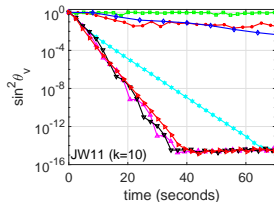
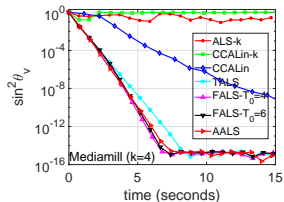
- Quality measures:

$$\sin^2 \theta_u \triangleq \sin^2 \theta_{\max}(\Phi_t, \mathbf{U}), \quad \sin^2 \theta_v \triangleq \sin^2 \theta_{\max}(\Psi_t, \mathbf{V})$$

Experiments



Experiments



Thank you!