

Improved Estimation of Canonical Vectors in Canonical Correlation Analysis

Nicholas Asendorf, Ph.D.

`asendorf@umich.edu`

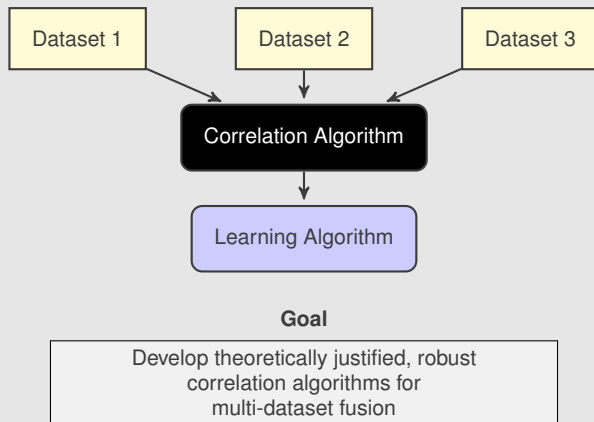
Prof. Raj Nadakuditi

`rajnrao@umich.edu`

Department of Electrical Engineering and Computer Science
University of Michigan

Asilomar Conference on Signals, Systems, and Computers

November 11, 2015



A Myriad of Applications

Multiple Datasets

- * Audio-Video
- * Audio-Audio

Machine Learning

- * emotion identification
- * shopping predictions
- * music genre classification

Medical Signal Processing

- * MRI, fMRI, EEG, MEG, etc.



A Myriad of Applications

Multiple Datasets

- * Audio-Video
- * Audio-Audio

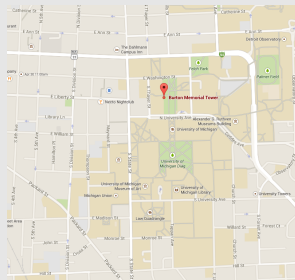
Machine Learning

- * emotion identification
- * shopping predictions
- * music genre classification



Medical Signal Processing

- * MRI, fMRI, EEG, MEG, etc.



A Myriad of Applications

Multiple Datasets

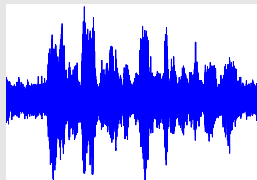
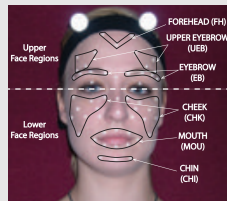
- * Audio-Video
- * Audio-Audio

Machine Learning

- * emotion identification
- * shopping predictions
- * music genre classification

Medical Signal Processing

- * MRI, fMRI, EEG, MEG, etc.



A Myriad of Applications

Multiple Datasets

- * Audio-Video
- * Audio-Audio

Machine Learning

- * emotion identification
- * shopping predictions
- * music genre classification

Medical Signal Processing

- * MRI, fMRI, EEG, MEG, etc.



Multiple Datasets

- * Audio-Video
- * Audio-Audio

Machine Learning

- * emotion identification
- * shopping predictions
- * music genre classification

Medical Signal Processing

- * MRI, fMRI, EEG, MEG, etc.



- * disco influences
- * danceable grooves
- * repetitive melodic phrasing
- * extensive vamping
- * minor key tonality

Canonical Correlation Analysis

What is it?

- * Dimensionality reduction algorithm for exactly 2 datasets, X, Y
- * Correlation coefficients, linear transformations

What is it not?

- * Data fusion algorithm

Covariance matrices

- * $R_{xx} = \mathbb{E} [xx^H]$
- * $R_{yy} = \mathbb{E} [yy^H]$
- * $R_{xy} = \mathbb{E} [xy^H]$

Variable Transformation

- * $f = R_{xx}^{1/2} w_x$
- * $g = R_{yy}^{1/2} w_y$

Optimization problem

$$\begin{array}{ll} \underset{w_x, w_y}{\operatorname{argmax}} & \rho = w_x^H R_{xy} w_y \\ \text{subject to} & w_x^H R_{xx} w_x = 1 \\ & w_y^H R_{yy} w_y = 1 \end{array}$$

Canonical Correlation Analysis

What is it?

- * Dimensionality reduction algorithm for exactly 2 datasets
- * Correlation coefficients, linear transformations

What is it not?

- * Data fusion algorithm

Optimization problem

$$\begin{aligned} \underset{f,g}{\operatorname{argmax}} \quad & \rho = f^H \underbrace{R_{xx}^{-1/2} R_{xy} R_{yy}^{-1/2}}_{C_{cca}} g \\ \text{subject to} \quad & \|f\|_2 = 1, \|g\|_2 = 1 \end{aligned}$$

Canonical Vectors

- * $w_x = R_{xx}^{-1/2} f$
- * $w_y = R_{yy}^{-1/2} g$

Insight

$$\# \text{ correlated signals} = k = \operatorname{rank}(C_{cca})$$

Training Datasets

- * $X = [x_1, \dots, x_n]$

- * $Y = [y_1, \dots, y_n]$

Sample Covariance Matrices

- * $\hat{R}_{xx} = \frac{1}{n}XX^H$

- * $\hat{R}_{yy} = \frac{1}{n}YY^H$

- * $\hat{R}_{xy} = \frac{1}{n}XY^H$

Estimate

$$\begin{aligned}\hat{C}_{cca} &= \hat{R}_{xx}^{-1/2} \hat{R}_{xy} \hat{R}_{yy}^{-1/2} \\ &= \hat{F} \hat{K} \hat{G}^H\end{aligned}$$

Data SVDs

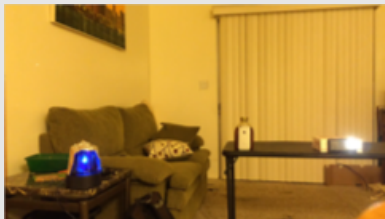
- * $X = \hat{U}_x \hat{\Sigma}_x \hat{V}_x^H$

- * $Y = \hat{U}_y \hat{\Sigma}_y \hat{V}_y^H$

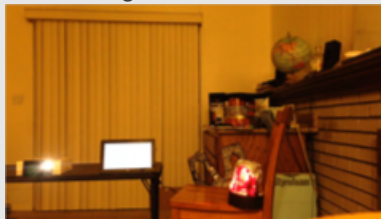
- * $\sigma(\hat{C}_{cca}) = \sigma(\hat{V}_x^H \hat{V}_y)$

Motivational Example - Flashing Light Video

Left Camera



Right Camera



System parameters

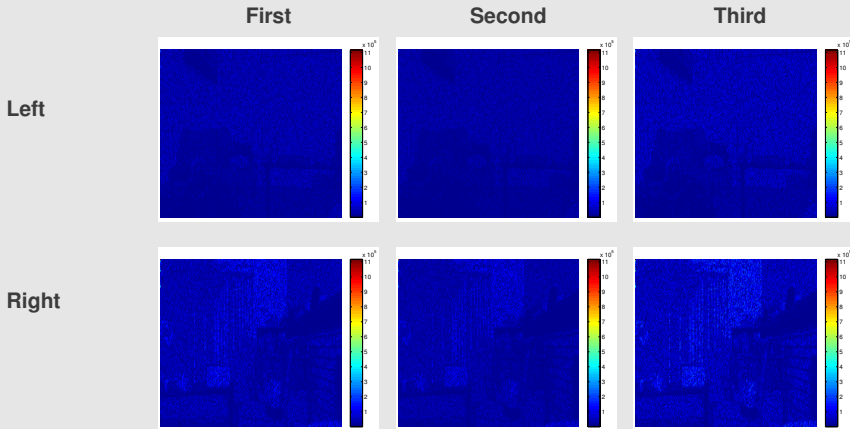
- * Vectorize 135×240 image $\Rightarrow p = q = 32400$ pixels
- * 30 fps @ 30 seconds $\Rightarrow n = 900$ frames

Goal

Identify correlated pixels between
camera views

Empirical CCA Results - Canonical Vectors

* After 900 frames = 30 seconds of video



Linear Subspace Model

$$x_i = U_x s_{x,i} + z_{x,i}$$

$$y_i = U_y s_{y,i} + z_{y,i}$$

Parameters

- * $U_x^H U_x = I_{k_x}, U_y^H U_y = I_{k_y}$
- * $z_{x,i} \stackrel{\text{i.i.d.}}{\sim} \mathcal{CN}(0, I_p), z_{y,i} \stackrel{\text{i.i.d.}}{\sim} \mathcal{CN}(0, I_q)$
- * $\mathbb{E} \left[\begin{bmatrix} s_{x,i} \\ s_{y,i} \end{bmatrix} \begin{bmatrix} s_{x,i}^H & s_{y,i}^H \end{bmatrix} \right] = \begin{bmatrix} \Theta_x & K_{xy} \\ K_{xy}^H & \Theta_y \end{bmatrix}$
- * $K_{xy} = \Theta_x^{1/2} P_{xy} \Theta_y^{1/2}$
- * $\Theta_x = \text{diag} \left(\left(\theta_1^{(x)} \right)^2, \dots, \left(\theta_{k_x}^{(x)} \right)^2 \right), \Theta_y = \text{diag} \left(\left(\theta_1^{(y)} \right)^2, \dots, \left(\theta_{k_y}^{(y)} \right)^2 \right)$
- * P_{xy} contains correlations ρ_{kj} between signals of x_i and y_i
- * $\tilde{K}_{xy} = (\Theta_x + I_{k_x})^{-1/2} K_{xy} (\Theta_y + I_{k_y})^{-1/2}$, with singular values $\kappa_1, \dots, \kappa_{\min(k_x, k_y)}$

Not all singular vectors are informative! (Nadakuditi, 2011)

- * $\sigma_i \left(\widehat{C}_{cca} \right) = \sigma_i \left(\widehat{V}_x^H \widehat{V}_y \right)$
- * Trim data SVD's to only use informative components

1. Trim data SVD's: $X = \widehat{U}_x \widehat{\Sigma}_x \widehat{V}_x^H$ and $Y = \widehat{U}_y \widehat{\Sigma}_y \widehat{V}_y^H$

- * $\mathring{U}_x = \widehat{U}_x \left(:, 1 : \widehat{k}_x \right), \mathring{U}_y = \widehat{U}_y \left(:, 1 : \widehat{k}_y \right)$

- * $\mathring{V}_x = \widehat{V}_x \left(:, 1 : \widehat{k}_x \right), \mathring{V}_y = \widehat{V}_y \left(:, 1 : \widehat{k}_y \right)$

2. Form $\widehat{C}_{icca} = \mathring{U}_x \mathring{V}_x^H \mathring{V}_y \mathring{U}_y$

3. Take SVD: $\widehat{C}_{icca} = \widetilde{F} \widetilde{K} \widetilde{G}^H$

4. $\widehat{\rho}_{icca}^{(i)} = \widetilde{k}_i$

5. $\widetilde{w}_x^{(i)} = \widehat{R}_{xx}^{-1/2} \widetilde{f}_i$

6. $\widetilde{w}_y^{(i)} = \widehat{R}_{yy}^{-1/2} \widetilde{g}_i$

Motivation

- * We expect ICCA to outperform CCA
- * However, we expect ICCA to be suboptimal because we substitute \hat{U}_x and \hat{U}_y without considering accuracy

Proposed Estimate

$$\hat{w}_{x,i}^{\text{icca+}} = \hat{U}_x \text{diag}(\lambda_{x,i}^{\text{opt}}) \hat{U}_{\tilde{K}}(:, i)$$

$$\hat{w}_{y,i}^{\text{icca+}} = \hat{U}_y \text{diag}(\lambda_{y,i}^{\text{opt}}) \hat{V}_{\tilde{K}}(:, i),$$

Optimization Problem

$$\lambda_{x,i}^{\text{opt}} = \underset{\lambda_x}{\operatorname{argmin}} \left\| w_x^{(i)} - \hat{U}_x \text{diag}(\lambda_x) \hat{U}_{\tilde{K}}(:, i) \right\|_F$$

$$\lambda_{y,i}^{\text{opt}} = \underset{\lambda_y}{\operatorname{argmin}} \left\| w_y^{(i)} - \hat{U}_y \text{diag}(\lambda_y) \hat{V}_{\tilde{K}}(:, i) \right\|_F.$$

Proposition

The solutions to the previous optimization problem is given by

$$\lambda_{x,i}^{opt} = \mathbf{diag} \left(\mathring{U}_x^H U_x (\Theta_x + I_{k_x})^{-1/2} \right)$$

$$\lambda_{y,i}^{opt} = \mathbf{diag} \left(\mathring{U}_y^H U_y (\Theta_y + I_{k_y})^{-1/2} \right).$$

Results from random matrix theory: Eigenvector accuracy

$$|\langle \hat{u}_i, u_i \rangle|^2 \xrightarrow{\text{a.s.}} \begin{cases} \frac{-2\varphi_{\mu_Z}(\rho)}{\theta_i^2 D'_{\mu_Z}(\rho)} & \theta_i^2 > 1/D_{\mu_Z}(b^+) \\ 0 & \text{otherwise} \end{cases}$$

$$D_{\mu}(z) =: \left[\int \frac{z}{z^2 - t^2} d_{\mu}(t) \right] \times \left[c_x \int \frac{z}{z^2 - t^2} d_{\mu}(t) + \frac{1 - c_x}{z} \right]$$

$$\varphi_{\mu}(z) =: \int \frac{z}{z^2 - t^2} d_{\mu}(t)$$

Theorem

The solution to the optimal ICCA+weights exhibits the following behavior in the asymptotic regime where $p, q, n \rightarrow \infty$ with $p/n \rightarrow c_x$ and $q/n \rightarrow c_y$.

For $i = 1, \dots, k_x$,

$$\lambda_{x,opt}^{(i)} \xrightarrow{a.s.} \begin{cases} D_{\mu_{Z_x}}(\sigma_x^{(i)}) \sqrt{\frac{-2\varphi_{\mu_{Z_x}}(\sigma_x^{(i)})}{D'_{\mu_{Z_x}}(\sigma_x^{(i)})(1+D_{\mu_{Z_x}}(\sigma_x^{(i)})}} & \text{if } (\theta_i^{(x)})^2 > 1/D_{\mu_{Z_x}}(b_x^+) \\ 0 & \text{otherwise} \end{cases}$$

where $\sigma_x^{(i)} = D_{\mu_{Z_x}}^{-1} \left(1 / (\theta_i^{(x)})^2 \right)$ and the D -transform. A similar expression for $\lambda_{y,opt}$ exists by replacing subscripts of x with y .

Estimation using training data

- * Using data X, Y , we can estimate $\widehat{D}, \widehat{D}', \widehat{\varphi}_{\mu}$, and $\widehat{\varphi}_{\tilde{\mu}}$

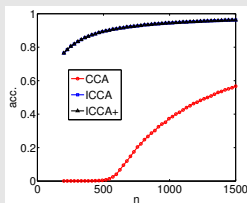
Data SVDs

- * $X = \hat{U}_x \hat{\Sigma}_x \hat{V}_x^H$ with trimmed versions $\mathring{U}_x, \mathring{\Sigma}_x, \mathring{V}_x$
- * $U_{\tilde{K}}$ left singular vectors of \tilde{K}_{xy}
- * $\hat{U}_{\tilde{K}}$ left singular vectors of \hat{C}_{cca}
- * $\mathring{U}_{\tilde{K}}$ left singular vectors of \hat{C}_{icca}

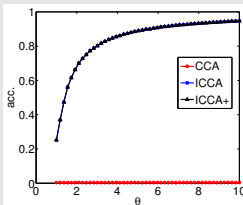
Estimate	
Population	$W_x = U_x (\Theta_x + I_{k_x})^{-1/2} U_{\tilde{K}}$
Empirical CCA	$\widehat{W}_x^{cca} = \hat{U}_x (\hat{\Sigma}_x)^{-1} \hat{U}_{\tilde{K}}$
ICCA	$\widehat{W}_x^{icca} = \mathring{U}_x \mathring{\Sigma}_x^{-1} \mathring{U}_{\tilde{K}}$
ICCA+	$\widehat{W}_x^{icca+} = \mathring{U}_x \Lambda_x^{\text{opt}} \mathring{U}_{\tilde{K}}$

Numerical Simulations

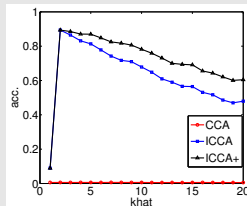
- * Rank-2 setting $k_x = k_y = 2$, $p = 200$, $q = 250$,
- * $\Theta_x = \Theta_y = \mathbf{diag}(16, 4)$, $P_{xy} = \mathbf{diag}(0.9, 0.5)$
- * Plot the accuracy of the first canonical vector for empirical CCA, ICCA, and ICCA+



(a) n sweep



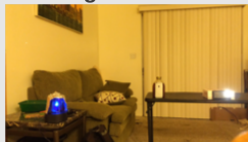
(b) θ sweep



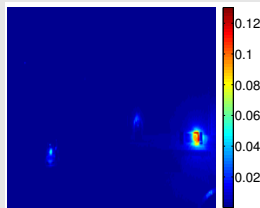
(c) \hat{k} sweep

Canonical Vector Accuracy Experiment

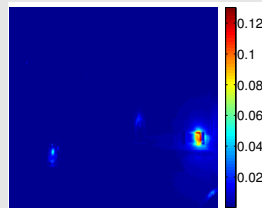
Original Scene



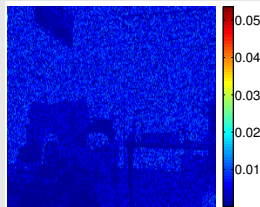
ICCA



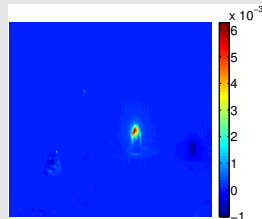
ICCA+



Empirical CCA



ICCA - ICCA+



- * Improved canonical vector accuracy in low-sample, low-SNR regime in Canonical Correlation Analysis (CCA)
- * Proposed a new algorithm: ICCA+
- * ICCA+ uses insights from random matrix theory to optimally scale sample eigenvectors setting
- * ICCA+ is more robust to over estimation of number of signals