

Learning From Data

Lecture 11: Unsupervised Learning III

Yang Li yangli@sz.tsinghua.edu.cn

TBSI

May 24, 2024

Today's Lecture

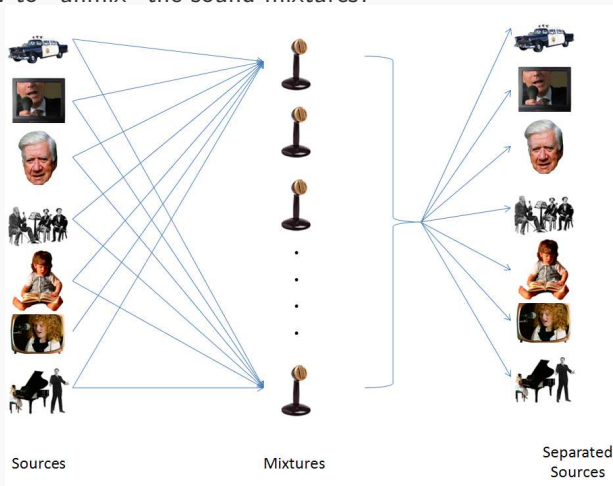
Unsupervised Learning (Part III)

- ▶ Independent Component Analysis (ICA)
- ▶ Canonical Correlation Analysis (CCA)

Independent Component Analysis

The cocktail party problem

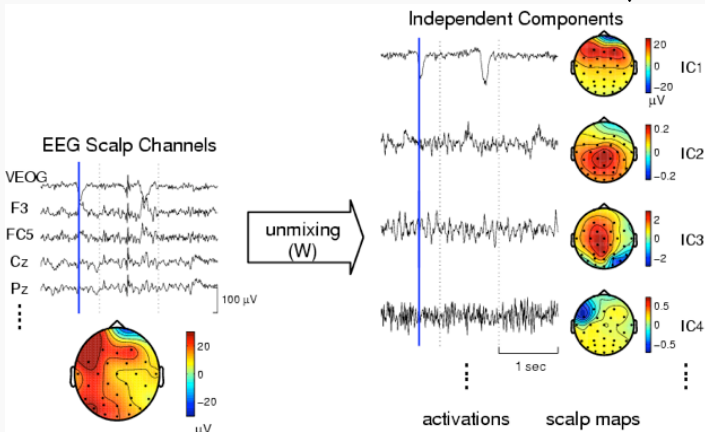
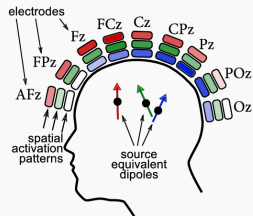
- ▶ n microphones at different locations of the room, each recording a mixture of n sound sources
- ▶ How to “unmix” the sound mixtures?



Sample audio: https://cnl.salk.edu/~tewon/Blind/blind_audio.html,
<http://www.kecl.ntt.co.jp/icl/signal/sawada/demo/bss2to4/index.html>

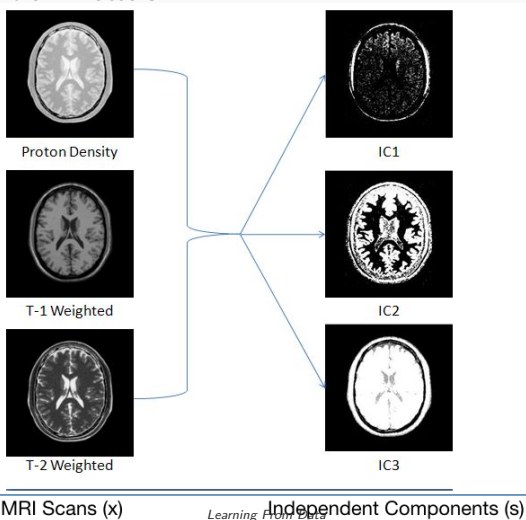
EEG Analysis

- ▶ Electrodes on patient scalp measure a mixture of different brain activations
- ▶ Finding independent activation sources helps removing artifacts in the signal



Brain imaging

- ▶ Different brain matters: gray matter, white matter, cerebrospinal fluid (CSF), fat, muscle/skin, glial matter etc.
- ▶ An MRI scan is a mixture of magnetic response signals from different brain matters



Problem Model

Case: $n = 2$

- ▶ Observed random variables: x_1, x_2
- ▶ Independent sources: $s_1, s_2 \in \mathbb{R}$

$$x_1 = a_{11}s_1 + a_{12}s_2$$

$$x_2 = a_{21}s_1 + a_{22}s_2$$

Problem Model

Case: $n = 2$

- ▶ Observed random variables: x_1, x_2
- ▶ Independent sources: $s_1, s_2 \in \mathbb{R}$

$$x_1 = a_{11}s_1 + a_{12}s_2$$

$$x_2 = a_{21}s_1 + a_{22}s_2$$

A is called the **mixing matrix**

$$x = As$$

Problem Model

Case: $n = 2$

- ▶ Observed random variables: x_1, x_2
- ▶ Independent sources: $s_1, s_2 \in \mathbb{R}$

$$x_1 = a_{11}s_1 + a_{12}s_2$$

$$x_2 = a_{21}s_1 + a_{22}s_2$$

A is called the **mixing matrix**

$$x = As$$

The blind source separation (cocktail party) problem

Given repeated observation $\{x^{(i)}; i = 1, \dots, m\}$, recover sources $s^{(i)}$ that generated the data ($x^{(i)} = As^{(i)}$)

Independent Component Analysis (ICA)

The blind source separation (cocktail party) problem

Given repeated observation $\{x^{(i)}; i = 1, \dots, m\}$, recover sources $s^{(i)}$ that generated the data ($x^{(i)} = As^{(i)}$)

Let $W = A^{-1}$ be the **unmixing matrix**

Goal of ICA: Find W , such that given $x^{(i)}$, the sources can be recovered by $s^{(i)} = Wx^{(i)}$

$$W = \begin{bmatrix} -w_1^T & - \\ \vdots & \\ -w_n^T & - \end{bmatrix}$$

Is W unique for a given set of observations ?

ICA Ambiguities

Assume data is **non Gaussian**, ICA has two ambiguities:

- ▶ Variance of the sources: *We can fix the magnitude of s_i by setting $\mathbb{E}[s_i^2] = 1$*

ICA Ambiguities

Assume data is **non Gaussian**, ICA has two ambiguities:

- ▶ Variance of the sources: *We can fix the magnitude of s_i by setting $\mathbb{E}[s_i^2] = 1$*
- ▶ Order of the sources s_1, \dots, s_n :
Let P be a permutation matrix, then we have $x = APP^{-1}s$.

ICA Ambiguities

Assume data is **non Gaussian**, ICA has two ambiguities:

- ▶ Variance of the sources: *We can fix the magnitude of s_i by setting $\mathbb{E}[s_i^2] = 1$*
- ▶ Order of the sources s_1, \dots, s_n :
Let P be a permutation matrix, then we have $x = APP^{-1}s$.

ICA Ambiguities

Assume data is **non Gaussian**, ICA has two ambiguities:

- ▶ Variance of the sources: *We can fix the magnitude of s_i by setting $\mathbb{E}[s_i^2] = 1$*
- ▶ Order of the sources s_1, \dots, s_n :
Let P be a permutation matrix, then we have $x = APP^{-1}s$.

Why is Gaussian data problematic?

ICA Ambiguities

Assume data is **non Gaussian**, ICA has two ambiguities:

- ▶ Variance of the sources: *We can fix the magnitude of s_j by setting $\mathbb{E}[s_j^2] = 1$*
- ▶ Order of the sources s_1, \dots, s_n :
Let P be a permutation matrix, then we have $x = APP^{-1}s$.

Why is Gaussian data problematic?

- ▶ The distribution of any rotation of Gaussian x has the same distribution as x .
- ▶ As long as at least one s_j is non-Gaussian, given enough data, we can recover the n independent sources.

Densities and Linear Transformations

Theorem 1

If random vector s has density p_s , and $x = As$ for a square, invertible matrix A , then the density of x is

$$p_x(x) = p_s(Wx) \cdot |W|$$

where $W = A^{-1}$.

ICA Algorithm

The joint distribution of *independent* sources $s = \{s_1, \dots, s_n\}$:

$$p(s) = \prod_{j=1}^n p_s(s_j)$$

ICA Algorithm

The joint distribution of *independent* sources $s = \{s_1, \dots, s_n\}$:

$$p(s) = \prod_{j=1}^n p_s(s_j)$$

The density of observation $x = As$ is:

$$p_x(x) = p_s(s)|W| = \prod_{j=1}^n p_s(s_j)|W| = \prod_{j=1}^n p_s(w_j^T x)|W|$$

ICA Algorithm

The joint distribution of *independent* sources $s = \{s_1, \dots, s_n\}$:

$$p(s) = \prod_{j=1}^n p_s(s_j)$$

The density of observation $x = As$ is:

$$p_x(x) = p_s(s) |W| = \prod_{j=1}^n p_s(s_j) |W| = \prod_{j=1}^n p_s(w_j^T x) |W|$$

Choose the sigmoid function $g(s) = \frac{1}{1+e^{-s}}$ as the *non-Gaussian* cdf for p_s , then

$$p_s(s) = g'(s)$$

ICA Algorithm

The joint distribution of *independent* sources $s = \{s_1, \dots, s_n\}$:

$$p(s) = \prod_{j=1}^n p_s(s_j)$$

The density of observation $x = As$ is:

$$p_x(x) = p_s(s) |W| = \prod_{j=1}^n p_s(s_j) |W| = \prod_{j=1}^n p_s(w_j^T x) |W|$$

Choose the sigmoid function $g(s) = \frac{1}{1+e^{-s}}$ as the *non-Gaussian* cdf for p_s , then

$$p_s(s) = g'(s)$$

This appears to be a heuristic choice, yet it can be justified rigorously in other interpretations.

ICA Algorithm

Given i.i.d. training samples $\{x^{(1)}, \dots, x^{(m)}\}$, the log likelihood is

$$\begin{aligned}l(W) &= \sum_{i=1}^m \log(p_x(x^{(i)})) = \sum_{i=1}^m \log\left(\prod_{j=1}^n p_s(w_j^T x) |W|\right) \\ &= \sum_{i=1}^m \left(\sum_{j=1}^n \log g'(w_j^T x^{(i)}) + \log |W| \right)\end{aligned}$$

ICA Algorithm

Given i.i.d. training samples $\{x^{(1)}, \dots, x^{(m)}\}$, the log likelihood is

$$\begin{aligned}l(W) &= \sum_{i=1}^m \log(p_x(x^{(i)})) = \sum_{i=1}^m \log\left(\prod_{j=1}^n p_s(w_j^T x) |W|\right) \\ &= \sum_{i=1}^m \left(\sum_{j=1}^n \log g'(w_j^T x^{(i)}) + \log |W| \right)\end{aligned}$$

Stochastic gradient ascent learning rule for sample $x^{(i)}$:

$$W := W + \alpha \left(\begin{bmatrix} 1 - 2g(w_1^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{bmatrix} x^{(i)T} + (W^T)^{-1} \right)$$

Check this at home!

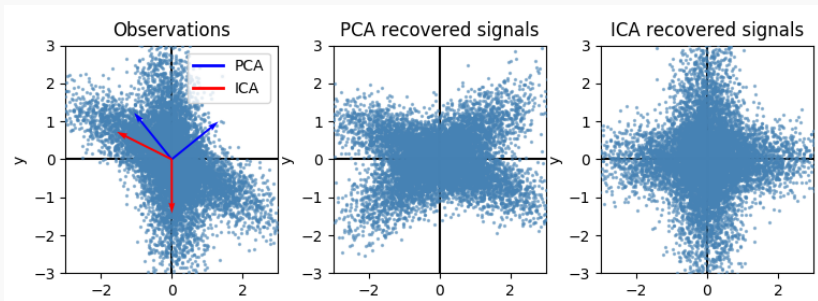
Theoretical Motivation of ICA

- ▶ Originally proposed by Jutten & Herault (1991) ¹*90 years later than PCA*
- ▶ Equivalent to learning projection directions w_1, \dots, w_n that
 - ▶ maximize the **sum of non-gaussianity** of the projected signals
 - ▶ minimize the **mutual information** of the projected signalsunder the constraint that $w_1^T x, \dots, w_n^T x$ are uncorrelated. ²

¹Christian Jutten, Jeanny Herault, Blind separation of sources, part I: An adaptive algorithm based on neuromimetic architecture, Signal Processing, Vol 24:1, 1991

²Hyvärinen, Aapo, and Erkki Oja. "Independent component analysis: algorithms and applications." Neural networks 13.4-5 (2000): 411-430.

ICA vs PCA



PCA	ICA
approximately Gaussian data	non-Gaussian data
removes correlation (low order dependence)	removes correlations and higher order dependence
ordered importance	all components are equally important
orthogonal	not orthogonal

Canonical Correlation Analysis

Canonical Correlation Analysis

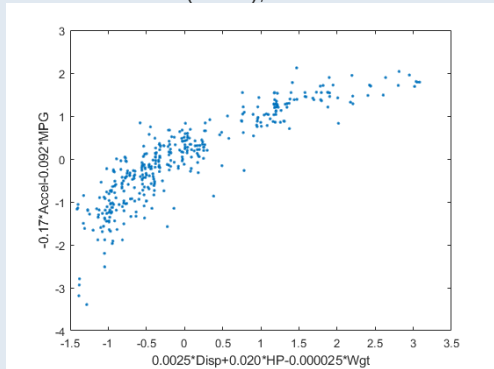
Canonical correlation analysis (CCA) finds the associations among two sets of variables.

Canonical Correlation Analysis

Canonical correlation analysis (CCA) finds the associations among two sets of variables.

Example: two sets of measurements of 406 cars:

- ▶ Specification: Engine displacement (Disp), horsepower (HP), weight (Wgt)
- ▶ Measurement: Acceleration (Accel), MPG



find important features that explain covariation between sets of variables

CCA Definitions

- ▶ Random vectors $X = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_1} \end{bmatrix}$ and $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{n_2} \end{bmatrix}$
- ▶ Covariance matrix $\Sigma_{XY} = \text{cov}(X, Y)$

CCA Definitions

- ▶ Random vectors $X = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_1} \end{bmatrix}$ and $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{n_2} \end{bmatrix}$
- ▶ Covariance matrix $\Sigma_{XY} = \text{cov}(X, Y)$
- ▶ CCA finds vectors a and b such that the random variables $a^T X$ and $b^T Y$ maximize the correlation

$$\rho = \text{corr}(a^T X, b^T Y)$$

CCA Definitions

- ▶ Random vectors $X = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_1} \end{bmatrix}$ and $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{n_2} \end{bmatrix}$
- ▶ Covariance matrix $\Sigma_{XY} = \text{cov}(X, Y)$
- ▶ CCA finds vectors a and b such that the random variables $a^T X$ and $b^T Y$ maximize the correlation

$$\rho = \text{corr}(a^T X, b^T Y)$$

- ▶ $U = a^T X$ and $V = b^T Y$ are called **the first pair of canonical variables**

CCA Definitions

- ▶ Random vectors $X = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_1} \end{bmatrix}$ and $Y = \begin{bmatrix} y_1 \\ \vdots \\ y_{n_2} \end{bmatrix}$
- ▶ Covariance matrix $\Sigma_{XY} = \text{cov}(X, Y)$
- ▶ CCA finds vectors a and b such that the random variables $a^T X$ and $b^T Y$ maximize the correlation

$$\rho = \text{corr}(a^T X, b^T Y)$$

- ▶ $U = a^T X$ and $V = b^T Y$ are called **the first pair of canonical variables**
- ▶ Subsequent pairs of canonical variables maximizes ρ while being *uncorrelated* with all previous pairs

Review: Singular Value Decomposition

A generalization of eigenvalue decomposition to rectangle ($m \times n$) matrices M .

$$M = U\Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T$$

- ▶ $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices
- ▶ $\Sigma \in \mathbb{R}^{m \times n}$ is a **rectangular diagonal matrix**.

Examples:

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \\ 0 & 0 & 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 \end{bmatrix}$$

Diagonal entries $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k$, $k = \min(n, m)$ are called **singular values of M** .

Review: Singular Value Decomposition

A non-negative real number σ is a singular value for $M \in \mathbb{R}^{m \times n}$ **if and only if** there exist unit-length $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$ such that

$$Mv = \sigma u$$

$$M^T u = \sigma v$$

u is called the **left singular vector** of σ , v is called the **right singular vector** of σ

Review: Singular Value Decomposition

A non-negative real number σ is a singular value for $M \in \mathbb{R}^{m \times n}$ **if and only if** there exist unit-length $u \in \mathbb{R}^m$ and $v \in \mathbb{R}^n$ such that

$$Mv = \sigma u$$

$$M^T u = \sigma v$$

u is called the **left singular vector** of σ , v is called the **right singular vector** of σ

Connection to eigenvalue decomposition

Given SVD of matrix $M = U\Sigma V^T$,

- ▶ $M^T M = (V\Sigma^T U^T)(U\Sigma V^T) = V(\Sigma^T \Sigma)V^T \leftarrow v_i$ is an eigenvector of $M^T M$ with eigenvalue σ_i^2
- ▶ $MM^T = (U\Sigma V^T)(V^T \Sigma^T U) = U(\Sigma \Sigma^T)U^T \leftarrow u_i$ is an eigenvector of MM^T with eigenvalue σ_i^2

CCA Derivations

The original problem:

$$(a_1, b_1) = \underset{a \in \mathbb{R}^{n_1}, b \in \mathbb{R}^{n_2}}{\operatorname{argmax}} \operatorname{corr}(a^T X, b^T Y) \quad (1)$$

CCA Derivations

The original problem:

$$(a_1, b_1) = \underset{a \in \mathbb{R}^{n_1}, b \in \mathbb{R}^{n_2}}{\operatorname{argmax}} \operatorname{corr}(a^T X, b^T Y) \quad (1)$$

Assume $\mathbb{E}[x_1] = \dots = \mathbb{E}[x_{n_1}] = \mathbb{E}[y_1] = \dots = \mathbb{E}[y_{n_2}] = 0$,

$$\begin{aligned} \operatorname{corr}(a^T X, b^T Y) &= \frac{\mathbb{E}[(a^T X)(b^T Y)^T]}{\sqrt{\mathbb{E}[(a^T X)^2]\mathbb{E}[(b^T Y)^2]}} \\ &= \frac{a^T \Sigma_{XY} b}{\sqrt{a^T \Sigma_{XX} a} \sqrt{b^T \Sigma_{YY} b}} \end{aligned}$$

CCA Derivations

The original problem:

$$(a_1, b_1) = \underset{a \in \mathbb{R}^{n_1}, b \in \mathbb{R}^{n_2}}{\operatorname{argmax}} \operatorname{corr}(a^T X, b^T Y) \quad (1)$$

Assume $\mathbb{E}[x_1] = \dots = \mathbb{E}[x_{n_1}] = \mathbb{E}[y_1] = \dots = \mathbb{E}[y_{n_2}] = 0$,

$$\begin{aligned} \operatorname{corr}(a^T X, b^T Y) &= \frac{\mathbb{E}[(a^T X)(b^T Y)^T]}{\sqrt{\mathbb{E}[(a^T X)^2]\mathbb{E}[(b^T Y)^2]}} \\ &= \frac{a^T \Sigma_{XY} b}{\sqrt{a^T \Sigma_{XX} a} \sqrt{b^T \Sigma_{YY} b}} \end{aligned}$$

(1) is equivalent to:

$$(a_1, b_1) = \underset{\substack{a \in \mathbb{R}^{n_1}, b \in \mathbb{R}^{n_2} \\ a^T \Sigma_{XX} a = b^T \Sigma_{YY} b = 1}}{\operatorname{argmax}} a^T \Sigma_{XY} b \quad (2)$$

CCA Derivations

CCA Derivations

Define $\Omega \in \mathbb{R}^{n_1 \times n_2}$, $c \in \mathbb{R}^{n_1}$ and $d \in \mathbb{R}^{n_2}$,

$$\Omega = \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$$

$$c = \Sigma_{XX}^{\frac{1}{2}} a$$

$$d = \Sigma_{YY}^{\frac{1}{2}} b$$

(2) can be written as

$$(c_1, d_1) = \underset{\substack{c \in \mathbb{R}^{n_1}, d \in \mathbb{R}^{n_2} \\ \|c\|^2 = \|d\|^2 = 1}}{\operatorname{argmax}} c^T \Omega d \quad (3)$$

CCA Derivations

Define $\Omega \in \mathbb{R}^{n_1 \times n_2}$, $c \in \mathbb{R}^{n_1}$ and $d \in \mathbb{R}^{n_2}$,

$$\Omega = \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$$

$$c = \Sigma_{XX}^{\frac{1}{2}} a$$

$$d = \Sigma_{YY}^{\frac{1}{2}} b$$

(2) can be written as

$$(c_1, d_1) = \underset{\substack{c \in \mathbb{R}^{n_1}, d \in \mathbb{R}^{n_2} \\ \|c\|^2 = \|d\|^2 = 1}}{\operatorname{argmax}} c^T \Omega d \quad (3)$$

(c_1, d_1) can be solved by SVD, then the first pair of canonical variables are

$$a_1 = \Sigma_{XX}^{-\frac{1}{2}} c_1, \quad b_1 = \Sigma_{YY}^{-\frac{1}{2}} d_1$$

CCA Derivations

$$(c_1, d_1) = \underset{\substack{c \in \mathbb{R}^{n_1}, d \in \mathbb{R}^{n_2} \\ \|c\|^2 = \|d\|^2 = 1}}{\operatorname{argmax}} c^T \Omega d$$

Proposition 1

c_1 and d_1 are the left and right unit singular vectors of Ω with the largest singular value.

CCA Derivations

$$(c_1, d_1) = \underset{\substack{c \in \mathbb{R}^{n_1}, d \in \mathbb{R}^{n_2} \\ \|c\|^2 = \|d\|^2 = 1}}{\operatorname{argmax}} c^T \Omega d$$

Proposition 1

c_1 and d_1 are the left and right unit singular vectors of Ω with the largest singular value.

Theorem 2

c_i and d_i are the left and right unit singular vectors of Ω with the i th largest singular value.

CCA Algorithm

Input: Covariance matrices for centered data X and Y :

- ▶ Σ_{XY} , invertible Σ_{XX} and Σ_{YY}
- ▶ Dimension $k \leq \min(n_1, n_2)$

Output: CCA projection matrices A_k and B_k :

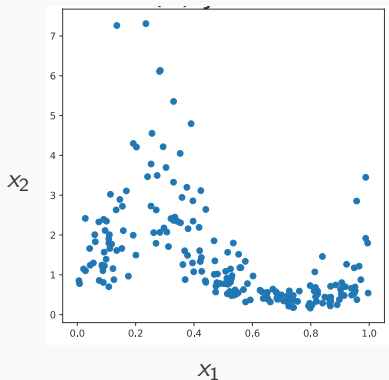
- ▶ Compute $\Omega = \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}$
- ▶ Compute SVD decomposition of Ω

$$\Omega = \begin{bmatrix} | & \cdots & | \\ c_1 & \cdots & c_{n_1} \\ | & \cdots & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \\ & & & 0 \end{bmatrix} \begin{bmatrix} -d_1^T \\ \vdots \\ -d_{n_2}^T \end{bmatrix}$$

- ▶ $A_k = \Sigma_{XX}^{-\frac{1}{2}} [c_1, \dots, c_k]$ and $B_k = \Sigma_{YY}^{-\frac{1}{2}} [d_1, \dots, d_k]$

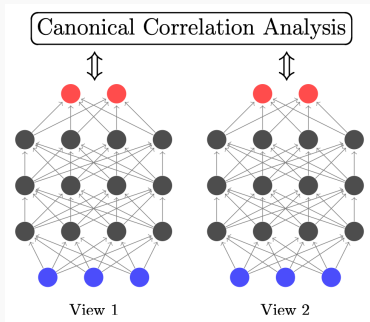
Discussion of CCA

- ▶ CCA only measures linear dependencies
- ▶ Non-linear generalizations:
 - ▶ Kernel CCA (KCCA)
 - ▶ Deep CCA (DCCA)
 - ▶ Maximal HGR Correlation



Non-linear dependency between x_1 and x_2

Deep Canonical Correlation Analysis (DCCA)



Andrew, Galen, et al. "Deep canonical correlation analysis." International conference on machine learning. PMLR, 2013.

Let $F_X = f(X; \theta_1)$, $G_Y = g(Y; \theta_2)$,

- ▶ Center features:

$$\bar{F}_X = F_X - \frac{1}{m} F_X^T \mathbf{1},$$

$$\bar{G}_Y = G_Y - \frac{1}{m} G_Y^T \mathbf{1}$$

- ▶ Define CCA Loss:

$$\theta_f^*, \theta_g^* = \operatorname{argmax}_{\theta_f, \theta_g} CCA(\bar{F}_X, \bar{G}_Y)$$

Maximize the total correlation of the top k components \implies

Maximize the sum of top k singular values of

$$\Omega = \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}}:$$

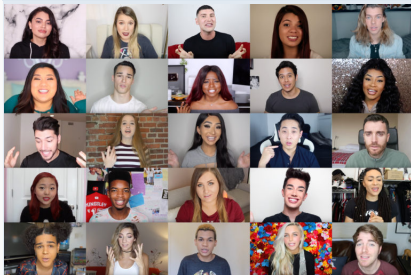
$$L_{CCA}(F_X, G_Y) = -\operatorname{tr}(\Omega^T \Omega)^{\frac{1}{2}}$$

- ▶ Update $\frac{\delta L_{CCA}(F_X, G_Y)}{\delta F_X}$, $\frac{\delta L_{CCA}(F_X, G_Y)}{\delta G_Y}$

Applications of CCA/DCCA

- ▶ Multiview clustering
Chaudhuri, Kamalika, et al. "Multi-view clustering via canonical correlation analysis." ICML 2009.
- ▶ Multimodal learning
Sun, Zhongkai, et al. "Learning relationships between text, audio, and video via deep canonical correlation for multimodal language analysis." AAAI 2020.

Multimodal sentiment analysis



Recognize speaker's emotion from videos using 3 modalities

- ▶ image
- ▶ text
- ▶ audio

(CMU-MOSEI dataset)

PCA, ICA and CCA

Linear Subspace Learning

Given high dimensional random vector \mathbf{x} , transform it to a low-dimensional vector \mathbf{y} through a projection matrix U :

$$\mathbf{y} = U^T \mathbf{x}$$

PCA, ICA and CCA

Linear Subspace Learning

Given high dimensional random vector \mathbf{x} , transform it to a low-dimensional vector \mathbf{y} through a projection matrix U :

$$y = U^T x$$

- ▶ PCA, ICA and CCA are all unsupervised linear subspace learning methods.

Name	What is U ?	goal	subspace
PCA	principal component (U)	remove (low order) correlation	single
ICA	unmixing matrix (W)	remove (high order) correlation	single
CCA	canonical projection matrices (A, B)	maximize correlation between feature pairs	paired