

Learning From Data

Lecture 12: Unsupervised Learning III

Yang Li yangli@sz.tsinghua.edu.cn

TBSI

December 16, 2022

Today's Lecture

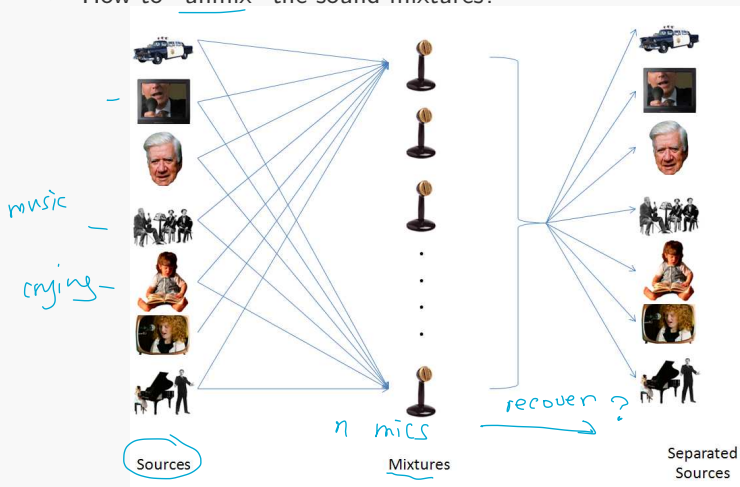
Unsupervised Learning (Part III)

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

Independent Component Analysis

The cocktail party problem

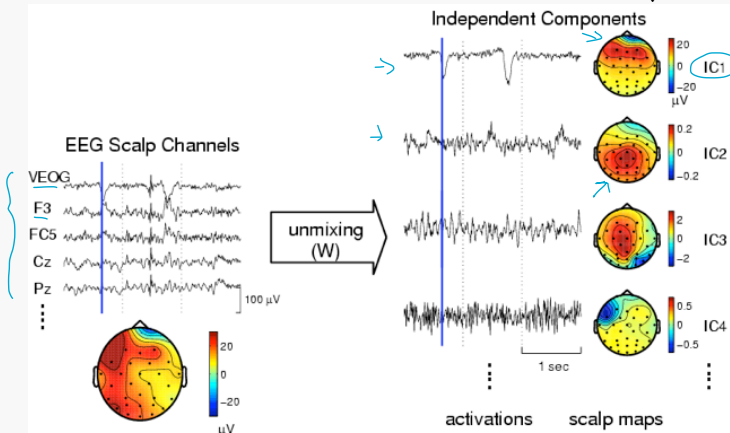
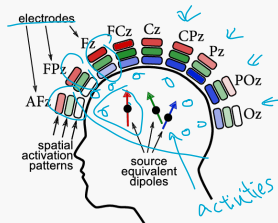
- ▶ n microphones at different locations of the room, each recording a mixture of n sound sources
- ▶ How to “unmix” the sound mixtures?



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

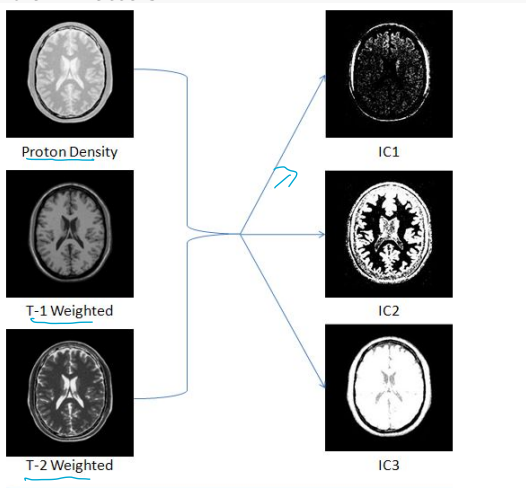
EEG Analysis

- ▶ 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

↑ # of independence source variables
→ # of observation variables

Case: $n = 2$

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

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

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

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** $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} s_1 \\ s_2 \end{pmatrix}$

$$x = As$$

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**

$$\begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \end{bmatrix} = x = As$$

$$s^{(i)} = \begin{bmatrix} s_1^{(i)} \\ s_2^{(i)} \end{bmatrix}$$

The blind source separation (cocktail party) problem

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

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 ($x^{(i)} = A s^{(i)}$)

\uparrow
mixing matrix

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

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

$$j \rightarrow \begin{bmatrix} s_j^{(i)} \\ \vdots \\ s_n^{(i)} \end{bmatrix} = \begin{bmatrix} -w_1^T & - \\ & \vdots \\ & -w_n^T & - \end{bmatrix} \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad W = \begin{bmatrix} -w_1^T & - \\ \vdots \\ -w_n^T & - \end{bmatrix}$$

$$s_j^{(i)} = w_j^T x^{(i)}$$

ICA Ambiguities

Question: Does assuming $\mathbb{E}[s_i^2] = 1$ resolve the scale ambiguity in $x = As$?
 A: no, because $-s_j$ gives the same result as s_j .

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

scale ► Variance of the sources: We can fix the magnitude of s_j by setting $\mathbb{E}[s_j^2] = 1$

$$\begin{aligned}
 x &= As \\
 x_j &= \sum_{i=1}^n a_{ji} s_i \quad \text{for all } j \\
 &= \sum_{i=1}^n a_{ji} \left(\frac{1}{c_j}\right) (c_j s_i) \quad \text{for any } c_j \neq 0 \\
 &\quad \underbrace{\hspace{1.5cm}}_{A'} \quad \underbrace{\hspace{1.5cm}}_{s'}
 \end{aligned}$$

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

ICA Ambiguities

$$\underbrace{\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}}_{P^{-1}} \begin{bmatrix} s_1 \\ s_2 \\ s_3 \end{bmatrix} = \begin{bmatrix} s_2 \\ s_3 \\ s_1 \end{bmatrix}$$

$\quad \quad \quad s \quad \quad \quad s'$

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 $\underline{s}_1, \dots, \underline{s}_n$:

Let P be a permutation matrix, then we have $x = APP^{-1}s$.

$n \times n$

$\underbrace{A}_{A'} \underbrace{P^{-1}s}_{s'}$

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?

independent sources $s = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$
 Given A , $x = As$

$x \sim \mathcal{N}(0, \underline{I})$
 $x \sim \mathcal{N}(0, AA^T)$
 rotation (orthogonal)

$\mathbb{E}[As] = A\mathbb{E}[s] = A \cdot 0 = 0$
 $\mathbb{E}[xx^T] = \mathbb{E}[(As)(As)^T] = \mathbb{E}[As s^T A^T] = A \underbrace{\mathbb{E}[s s^T]}_{\underline{I}} A^T = AA^T$

$x' = (A \cdot R)s$ $R^{-1} = R^T$

$x' \sim \mathcal{N}(0, \underbrace{(AR)(AR)^T}_{ARR^T A^T}) = \mathcal{N}(0, AA^T)$
 $\underbrace{ARR^T}_{\underline{I}} A^T = AA^T$

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?

- ▶ 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.

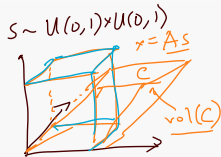
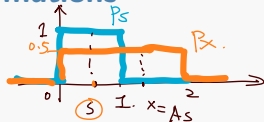
Densities and Linear Transformations

Let $s \sim \text{Uniform}(0, 1)$

Let $A = 2$

$x = A s \sim \text{Uniform}(0, 2)$

$$p_s(s) = 1. \quad \left\{ \begin{array}{l} p_x(x) = p_x(As) = 0.5 \\ p_s(s) \cdot |A^{-1}| = 1 \cdot \frac{1}{2} = 0.5 \end{array} \right.$$



$$p_x = \frac{1}{\text{vol}(C)} \quad \{x \in C\}$$

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

$$p_x(x) = p_s(Wx) \underbrace{|W|}_{\text{determinant}}$$

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

ICA Algorithm

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

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

ICA Algorithm

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

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

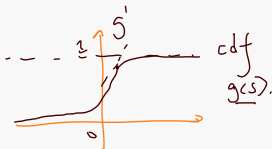
The density of observation $x = As$ is: $s_j = w_j^T x$.

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

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

$$\underline{p_s(s)} = \underline{g'(s)}$$

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

$$\underline{p_s(s)} = g'(s)$$

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

ICA Algorithm

pdf of s : $g'(s) = g(s)(1-g(s))$
 $\frac{d}{ds} g(s)$

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

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

$g(w_j^T x^{(i)})$ $(+ g(w_j^T x^{(i)}))$

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

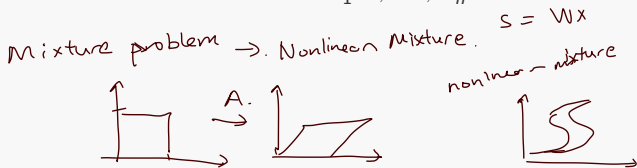
$$W := W + \alpha \left(\begin{bmatrix} 1 - 2g(w_1^T x^{(i)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(i)}) \end{bmatrix} x^{(i)T} + (W^T)^{-1} \right)$$

$\nabla W.$

Check this at home!

Theoretical Motivation of ICA

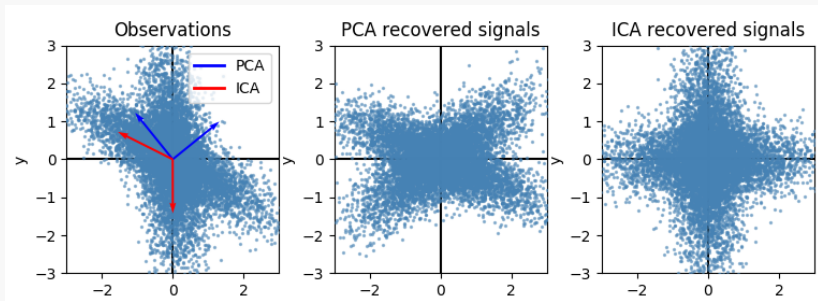
- Originally proposed by Jutten & Herault (1991) ¹ 90 years later than PCA
- Equivalent to learning projection directions $\underline{w}_1, \dots, \underline{w}_n$ that
 - maximize the sum of non-gaussianity of the projected signals
 - minimize the mutual information of the projected signals \rightarrow independent sources
 under 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

approximately Gaussian data

removes correlation (low order dependence)

ordered importance

orthogonal

ICA

non-Gaussian data

removes correlations and higher order dependence $E[x^4]$

all components are equally important

not orthogonal

Canonical Correlation Analysis

Canonical Correlation Analysis

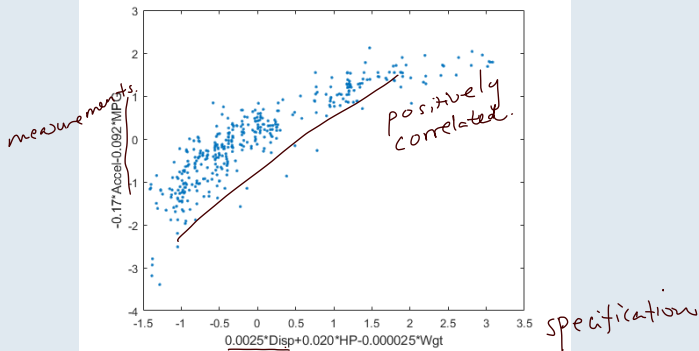
Canonical correlation analysis (CCA) finds the associations among two sets of variables.

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 $\underline{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_{n_1} \end{bmatrix}$ and $\underline{Y} = \begin{bmatrix} y_1 \\ \vdots \\ y_{n_2} \end{bmatrix}$

$\text{cov}(X)$ $\text{cov}(Y)$
 $\text{cov}(X, X)$ $= \Sigma_{YY}$
 Σ_{XX}

▶ Covariance matrix $\Sigma_{XY} = \underline{\text{cov}(X, Y)}$

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

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

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

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

- ▶ $U = a^T X$ 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

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

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

st. ...

- ▶ $U = a^T X$ 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

$$a_2 \perp a_1$$

Review: Singular Value Decomposition

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

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

Handwritten annotations: $(m \times n)$, $(m \times m)$, $(m \times n)$, $(n \times n)$, $r = \text{rank}$, $u_1 \dots u_m$, $v_1 \dots v_n$.

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

Examples:

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

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

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

parallel to eigenvector function
 $Au = \lambda u$

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

$$Mv = \sigma u$$

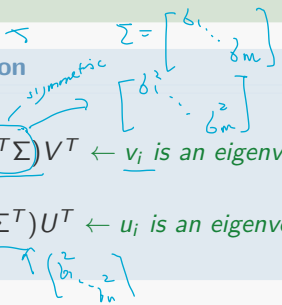
$$M^T u = \sigma v$$

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:

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

□

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

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

$\leftarrow a^T \mathbb{E}[XY^T] b$

□

CCA Derivations

The original problem:

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

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

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

(1) is equivalent to:

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

$$\underbrace{a^T \Sigma_{XX} a = b^T \Sigma_{YY} b = 1}$$

CCA Derivations

$$\max_{\substack{a \in \mathbb{R}^{n_1} \\ b \in \mathbb{R}^{n_2}}} a^T \Sigma_{xy} b. \quad (1)$$

$$\text{st. } \frac{a^T \Sigma_{xx} a = 1}{b^T \Sigma_{yy} b = 1}.$$

$$\begin{aligned} a^T \Sigma_{xy} b &= a^T \underbrace{\Sigma_{xx}^{-\frac{1}{2}} \Sigma_{xx}^{\frac{1}{2}}}_{\mathbf{I}} \Sigma_{xy} \Sigma_{yy}^{-\frac{1}{2}} \Sigma_{yy}^{\frac{1}{2}} b \\ &= \left((\Sigma_{xx}^{-\frac{1}{2}})^T a \right)^T \Sigma_{xy} \Sigma_{yy}^{-\frac{1}{2}} \Sigma_{yy}^{\frac{1}{2}} b \\ &= \underbrace{(\Sigma_{xx}^{-\frac{1}{2}} a)^T}_{\mathbf{c}^T} \underbrace{(\Sigma_{xx}^{-\frac{1}{2}} \Sigma_{xy} \Sigma_{yy}^{-\frac{1}{2}})}_{\mathbf{\Omega}} \underbrace{(\Sigma_{yy}^{\frac{1}{2}} b)}_{\mathbf{d}} \end{aligned}$$

$$a^T \Sigma_{xx} a = a^T \Sigma_{xx}^{\frac{1}{2}} \Sigma_{xx}^{\frac{1}{2}} a = \mathbf{c}^T \mathbf{c} = \|\mathbf{c}\|^2 = 1$$

$$b^T \Sigma_{yy} b = \mathbf{d}^T \mathbf{d} = \|\mathbf{d}\|^2 = 1$$

Therefore, (1) is equivalent to

$$\begin{aligned} \max_{\substack{c \in \mathbb{R}^{n_1} \\ d \in \mathbb{R}^{n_2}}} c^T \mathbf{\Omega} d \\ \text{st. } \|\mathbf{c}\|^2 = 1 \\ \|\mathbf{d}\|^2 = 1 \end{aligned}$$

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

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

CCA Derivations

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

$$\begin{aligned}\Omega &= \Sigma_{XX}^{-\frac{1}{2}} \Sigma_{XY} \Sigma_{YY}^{-\frac{1}{2}} \\ c &= \Sigma_{XX}^{\frac{1}{2}} a \Rightarrow a = \Sigma_{XX}^{-\frac{1}{2}} c \\ d &= \Sigma_{YY}^{\frac{1}{2}} b\end{aligned}$$

(2) can be written as

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

(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

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

Proposition 1

c_1 and d_1 are the left and right unit singular vectors of $\underline{\Omega}$ with the largest singular value.

Theorem 2

\underline{c}_i and \underline{d}_i are the left and right unit singular vectors of $\underline{\Omega}$ with the \underline{i} th largest singular value.

$$\underline{\Omega} = \underbrace{U} \Sigma \underbrace{V^T}$$

CCA Algorithm

Input: Covariance matrices for centered data X and Y :

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

Output: CCA projection matrices A_k and B_k :

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

$$\underline{\Omega} = \underbrace{\begin{bmatrix} | & \dots & | \\ c_1 & \dots & c_{n_1} \\ | & \dots & | \end{bmatrix}}_k \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \\ & & & 0 \end{bmatrix} \begin{bmatrix} -d_1^T \\ \vdots \\ -d_{n_2}^T \end{bmatrix} \left. \vphantom{\begin{bmatrix} -d_1^T \\ \vdots \\ -d_{n_2}^T \end{bmatrix}} \right\}^k$$

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

Discussion of CCA

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

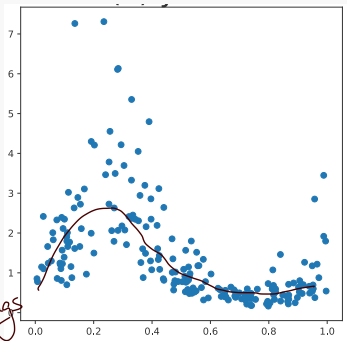
$$\begin{aligned} x &\rightarrow f(x) \\ y &\rightarrow g(y) \end{aligned}$$

$$\begin{cases} \max_{f, g} \mathbb{E}[f(x)^T g(y)] \\ \text{s.t. } \mathbb{E}[f] = \mathbb{E}[g] = 0; \mathbb{E}[f^2] = \mathbb{E}[g^2] = 1. \end{cases}$$

non-parametric mappings

$\phi(x_2)$

x_2



$x_1 \rightarrow \phi(x_1)$

Non-linear dependency between x_1 and x_2

x_1, y are discrete
 \rightarrow Alternating Directional Conditional Expectation (ADCE).

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 :

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

PCA, ICA and CCA

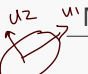


Linear Subspace Learning

$$\left. \begin{array}{l} \underline{x} \\ \underline{z} \end{array} \right\}$$

Given high dimensional random vector \underline{x} , transform it to a low-dimensional vector \underline{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	<u>principal component</u> (U)	remove (low order) correlation	<u>single</u>
 ICA	<u>unmixing matrix</u> (W)	remove (<u>high order</u>) correlation	single
 CCA	canonical projection matrices (A, B)	maximize correlation between feature pairs	<u>paired</u>