

# Canonical Correlation Analysis

## (1) Back ground:

Recall: Regression analysis concerns with relationship between a single response and a set of predictors.

What if we have more than one responses?

i) Multivariate Regression:  $Y_k = \sum_{i=1}^p \beta_{ki} X_i + \varepsilon_k, 1 \leq k \leq 2$

ii) Apply PCA on  $Y_k$ .      iii) CCA.

e.g. Variables related to arithmetic power and  
Variables related to reading power.

⇒ The advantage of CCA is that:

It seeks to identify and quantify linear associations between two sets of variables by find L.F. of variables maximally correlated, extract the valuable information about correlation.

## (2) Population CCA:

① Data: 1<sup>st</sup> group  $X^{(1)} = (X_{11} \dots X_{1p})^T$       suppose  $p \leq q$ .

2<sup>nd</sup> group  $X^{(2)} = (X_{21} \dots X_{2q})^T$

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}, E(X) = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \text{Var}(X) = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Prmk:  $p, q$  elements  $\Sigma_{12}$  measures the association between two sets. Commonly, we require two variables are homogeneous

Set  $U = a^T X^{(1)}$ ,  $V = b^T X^{(2)}$ . Then we obtain:

$$\text{Var}(U) = a^T \Sigma_{11} a, \quad \text{Var}(V) = b^T \Sigma_{22} b.$$

$$\text{Cov}(U, V) = a^T \Sigma_{12} b. \quad \text{We seek } a, b \text{ s.t.}$$

$$\text{maximizes: } \text{Corr}(U, V) = \frac{(a^T \Sigma_{12} b)^2}{(a^T \Sigma_{11} a)(b^T \Sigma_{22} b)} \quad (*)$$

i) The 1<sup>st</sup> pair  $(U_1, V_1)$  maximizes  $(*)$ , which's constrained by:  $\text{Cov}(U_1) = \text{Cov}(V_1) = 1$ .

ii) The 2<sup>nd</sup> pair  $(U_2, V_2)$  need to extract the max information of correlation which are uncorrelated with  $(U_1, V_1)$ , i.e. maximize  $(*)$ , and s.t.  
 $\text{Cov}(U_2) = \text{Cov}(V_2) = 1, \quad \text{Cov}(U_2, U_1) = \text{Cov}(U_2, V_1)$   
 $= \text{Cov}(V_2, U_1) = \text{Cov}(V_2, V_1) = 0.$

iii) The  $k^{\text{th}}$  pair  $(U_k, V_k)$  maximizes  $(*)$ , s.t.  
 have unit Var. uncorrelated with the first  $k-1$  pairs, maximizes  $(*)$ .

Prmk: Compare to PCA:

CCA	PCA
two sets	one set
2 proj.	1 proj.
max $\text{Corr}^2$	max Var



Thm. Denote  $\text{Corr}^2(u_k, v_k) = \rho_k^{*2}$ .  $\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$ .

Then  $a_k^T = e_k^T \Sigma_{11}^{-\frac{1}{2}}$ ,  $b_k^T = f_k^T \Sigma_{22}^{-\frac{1}{2}}$  where  $(e_k^T, f_k)$  is eigen-pair of  $\Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-\frac{1}{2}}$ .  $(e_k^T, f_k)$  is eigen-pair of  $\Sigma_{22}^{-\frac{1}{2}} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$ .

pf: 1) By Schwartz Ineqn. fix  $b$ :

$$\text{Corr}(a^T X_{(1)}, b^T X_{(2)}) \leq \frac{a^T \Sigma_{11} a (b^T \Sigma_{12} \Sigma_{11}^{-1} \Sigma_{22} b)}{(a^T \Sigma_{11} a) (b^T \Sigma_{22} b)}$$

$$= \frac{b^T \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} b}{b^T \Sigma_{22} b} \leq \rho_1^{*2} \text{ which is maximum}$$

eigenvalue of  $\Sigma_{22}^{-\frac{1}{2}} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$ .

Besides  $\tilde{a}_1 = c \Sigma_{11}^{-1} \Sigma_{12} \tilde{b}_1$ . " $=$ " holds.

Normalized:  $a_1 = \Sigma_{11}^{-\frac{1}{2}} e_1 = \frac{1}{\sqrt{\lambda_1}} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} f_1$

$$b_1 = \Sigma_{22}^{-\frac{1}{2}} f_1, \quad \lambda_1 = \rho_1^{*2}.$$

2) By induction: suppose  $1 \leq k \leq n-1$  holds.

$$\text{Set } u_k = \Sigma_{11}^{-\frac{1}{2}} b_k, \quad 1 \leq k \leq p, \quad u = \Sigma_{11}^{-\frac{1}{2}} b.$$

$$\text{From above: } \rho^2(u_n, v_n) \leq \tilde{u}^T T^T T \tilde{u}, \quad \tilde{u}^T \tilde{u} = 1,$$

$$\text{and } T = \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}, \quad \tilde{u} \in \text{span}^\perp(u_1, \dots, u_{n-1}).$$

$$\Rightarrow \tilde{u} = \sum_n^p c_k u_k = \sum_n^p c_k f_k, \quad \sum c_k^2 = 1.$$

$$\therefore \rho^2(u_n, v_n) \leq \sum_n^p c_k^2 \lambda_k \leq \lambda_n = \rho_n^{*2}.$$

where " $=$ " holds when  $\tilde{u} = f_n = u_n$ .

Rmk: i) Non-zero eigenvalues of  $T^T T$  is identical

with  $T T^T$ . Besides,  $\rho_k = \frac{1}{\sqrt{\lambda_k}} \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} f_k$ .

$$\lambda_k = \rho_k^{*2}.$$



ii) Note the matrix is symmetric. So  $e_i \perp e_j$ ,  $f_i \perp f_j$ ,  $i \neq j$ . But  $b_i \nperp b_j$  and  $a_i \nperp a_j$  may hold!

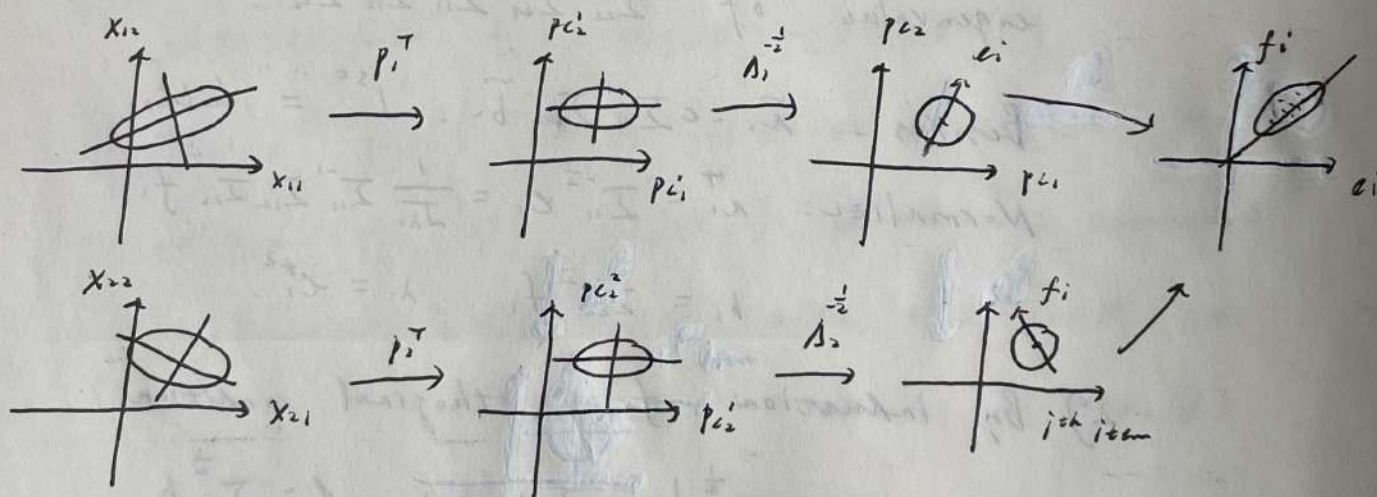
### Geometric Interpretation:

Let  $A = (a_1, \dots, a_p)_{p \times p}$ ,  $B = (b_1, \dots, b_p)_{2 \times p}$

$\Rightarrow U = A^T X_{(1)}$ ,  $V = B^T X_{(2)}$ ,  $A^T = E^T \Sigma_{11}^{-\frac{1}{2}}$ ,  $E = (e_1, \dots, e_p)$

Decompose  $\Sigma_{11}^{-\frac{1}{2}} = P_1 \Lambda_1^{-\frac{1}{2}} P_1^T$ ,  $\therefore U = E^T P_1 \Lambda_1^{-\frac{1}{2}} P_1^T X_{(1)}$

Note that  $P_1^T X_{(1)}$  is PC (analogous to  $X_{(2)}$ )



Rmk: The last  $E^T P_i$  is just rotation to select a direction to projection, which guarantee large correlation.

The direction may do nothing with PCs.

### ② Properties of Canonical Variables:

i)  $\text{COV}(U, X_{(1)}) = A^T \Sigma_{11}$ ,  $\text{COV}(V, X_{(2)}) = B^T \Sigma_{22}$

$\text{COV}(U, X_{(2)}) = A^T \Sigma_{12}$ ,  $\text{COV}(V, X_{(1)}) = B^T \Sigma_{21}$ .

$$\text{Corr}(U, X^{(1)}) = \text{Cov}(U, V_{11}^{-\frac{1}{2}} X_{(1)}) = A^T \Sigma_{11} V_{11}^{-\frac{1}{2}}$$

$$\text{Corr}(V, X^{(2)}) = \text{Cov}(V, V_{22}^{-\frac{1}{2}} X_{(2)}) = B^T \Sigma_{22} V_{22}^{-\frac{1}{2}}$$

$$\text{where } V_{11} = \text{diag } \Sigma_{11}, \quad V_{22} = \text{diag } \Sigma_{22}.$$

## ii) Invariant:

Consider model:  $X_{(1)}^* = P X_{(1)} + c_1$ ,  $X_{(2)}^* = Q X_{(2)} + c_2$ .

Then CCA on  $(X_{(1)}, X_{(2)})$  is essentially same to  $(X_{(1)}^*, X_{(2)}^*)$ . and  $a_i^* = P^T a_i$ ,  $b_i^* = Q^T b_i$ ,  $1 \leq i \leq p$ .

## iii) To ease the computation burden:

$$\text{Calculate: } \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} a = c^* a, \quad \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} b = c^* b.$$

$$\text{i.e. } |\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} - c^* I| = 0.$$

## iv) Some Interpretation:

- Compute correlation between canonical variables and original variables can be a way to determine the relative important of origin and canonical

- Canonical correlation generalizes the correlation between 2 variables to 2 groups variables.

$$|\text{Corr}^2(X_i^{(1)}, X_k^{(2)})| = |\text{Corr}^2(e_i^T X^{(1)}, e_k^T X^{(2)})| \leq c_i^{*2}$$

- In multiple correlation coefficient interpretation:

$c_k^{*2}$  is the proportion of Var of  $U_k = a_k^T X^{(1)}$  explained by the linear combination of  $X^{(2)}$ .

It's also explain  $V_k$  by  $X^{(1)}$ .



### ③ CCA for Standardization:

Consider  $Z^{(1)} = (Z_1^{(1)}, \dots, Z_p^{(1)})^T$ ,  $Z^{(2)} = (Z_1^{(2)}, \dots, Z_q^{(2)})^T$

$$\Rightarrow \begin{cases} U_k = A_k^T Z^{(1)} = e_k^T e_{11}^{-\frac{1}{2}} Z^{(1)} \\ V_k = b_k^T Z^{(2)} = f_k^T e_{22}^{-\frac{1}{2}} Z^{(2)} \end{cases} \text{ where}$$

$(e_k, e_k^*)$ ,  $(f_k, f_k^*)$  are eigenpairs of  $e_{11}^{-\frac{1}{2}} e_{12}^{-1} e_{22}^{-\frac{1}{2}} e_{21}$ ,  
 $e_{22}^{-\frac{1}{2}} e_{21}^{-1} e_{11}^{-\frac{1}{2}} e_{12}$ .  $e_k = \frac{1}{\sqrt{\lambda_k}} \cdot e_{11}^{-\frac{1}{2}} e_{12}^{-1} e_{22}^{-\frac{1}{2}} f_k$ .

Hint: CCA is unchanged under standardization.

i.e.  $\Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-\frac{1}{2}}$  and  $e_{11}^{-\frac{1}{2}} e_{12}^{-1} e_{22}^{-\frac{1}{2}} e_{21} e_{11}^{-\frac{1}{2}}$  will have same eigenvalues. Suppose  $e_k, e_k^*$  are correspond eigenvectors.  $a_k = \Sigma_{11}^{-\frac{1}{2}} e_k$ ,  $a_k^* = e_{11}^{-\frac{1}{2}} e_k^*$

$$\Rightarrow a_k^* = V_{11}^{-\frac{1}{2}} a_k. \quad V_{11} = \text{diag } \Sigma_{11}.$$

### (3) CCA on Samples:

Data:  $X = (X^{(1)} \mid X^{(2)}) =$

$$= \begin{pmatrix} X_1^{(1)T} & X_1^{(2)T} \\ \vdots & \vdots \\ X_n^{(1)T} & X_n^{(2)T} \end{pmatrix}$$

$$\begin{matrix} & \text{Var}_1^{(1)} & \text{Var}_2^{(1)} & \dots & \text{Var}_p^{(1)} & \text{Var}_1^{(2)} & \dots & \text{Var}_q^{(2)} \\ \text{item} \rightarrow & X_{11}^{(1)} & X_{12}^{(1)} & \dots & X_{1p}^{(1)} & X_{11}^{(2)} & \dots & X_{1q}^{(2)} \\ & \vdots & \vdots & & \vdots & \vdots & & \vdots \\ & X_{n1}^{(1)} & \dots & X_{np}^{(1)} & X_{n1}^{(2)} & \dots & X_{nq}^{(2)} \end{matrix}$$

Find CCA: i) Replace population list by empirical list.

ii) Replace  $\Sigma$  by  $S$ ,  $e$  by  $R$ .

### ① Matrixs of Error Approx:

$$\hat{U} = \hat{A}^T X^{(1)}, \quad \hat{V} = \hat{B}^T X^{(2)} \Rightarrow X^{(1)} = (\hat{A}^T)^{-1} \hat{U}, \quad X^{(2)} = (\hat{B}^T)^{-1} \hat{V}.$$

Note that:  $\text{COV}(\hat{U}, \hat{V}) = \begin{pmatrix} \hat{e}_1^* & & \\ & \ddots & \\ & & \hat{e}_p^* \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} = \hat{A}^T S_{12} \hat{B}$

Denote:  $(\hat{A}^T)^{-1} = (\hat{a}^{(1)} \dots \hat{a}^{(p)})$

$$(\hat{B}^T)^{-1} = (\hat{b}^{(1)} \dots \hat{b}^{(r)})$$

$$\Rightarrow S_{12} = \frac{1}{P} \sum_i \hat{e}_i^* \hat{a}^{(i)} \hat{b}^{(i)T}. \quad \text{Similarly, } \text{COV}(\hat{U}) = \text{COV}(\hat{V}) = I$$

$$\Rightarrow S_{11} = \frac{1}{P} \sum_i \hat{a}^{(i)} \hat{a}^{(i)T}, \quad S_{22} = \frac{1}{P} \sum_i \hat{b}^{(i)} \hat{b}^{(i)T}.$$

Rmk:  $X^{(1)} = (\hat{A}^T)^{-1} \hat{U} = \sum_{i=1}^p \hat{U}_i \hat{a}^{(i)} \Rightarrow \text{COV}(X^{(1)}, \hat{U}_i) = \hat{a}^{(i)}$

$\Rightarrow$  The first  $r$  columns contain samples COV of  $\hat{U}_1 \dots \hat{U}_r$  and  $X_1^{(1)} \dots X_p^{(1)}$ .

Similar for  $(\hat{B}^T)^{-1}$  and  $\hat{V}_1 \dots \hat{V}_r$ .

If only the first  $r$  canonical pairs are used:

$$\bar{X}^{(1)} = (\hat{a}^{(1)} \dots \hat{a}^{(r)}) \begin{pmatrix} \hat{U}_1 \\ \vdots \\ \hat{U}_r \end{pmatrix} \in \mathbb{R}^P \Rightarrow S_{12} \text{ is approxi}$$

$$\bar{X}^{(2)} = (\hat{b}^{(1)} \dots \hat{b}^{(r)}) \begin{pmatrix} \hat{V}_1 \\ \vdots \\ \hat{V}_r \end{pmatrix} \in \mathbb{R}^P \quad \text{by } \text{COV}(\bar{X}^{(1)}, \bar{X}^{(2)})$$

Residuals:

$$\begin{cases} S_{11} - \sum_{i=1}^r \hat{a}^{(i)} \hat{a}^{(i)T} = \sum_{i=r+1}^p \hat{a}^{(i)} \hat{a}^{(i)T} \\ S_{22} - \sum_{i=1}^r \hat{b}^{(i)} \hat{b}^{(i)T} = \sum_{i=r+1}^p \hat{b}^{(i)} \hat{b}^{(i)T} \\ S_{12} - \sum_{i=1}^r \hat{e}_i^* \hat{a}^{(i)} \hat{b}^{(i)T} = \sum_{i=r+1}^p \hat{e}_i^* \hat{a}^{(i)} \hat{b}^{(i)T} \end{cases}$$

Rmk: i) Large entries of residual matrix indicate a poor fit for correspond variables.

ii)  $S_{12}$  is better fit than  $S_{11}, S_{22}$ . Since

We can select  $r$  s.t.  $\hat{e}_k^*$  is small.



when  $k \geq r+1$ . But for  $S_{11}, S_{22}$ , the approxi. may not as good as  $S_{12}$ . since it can't be controlled

## ② Proportions of Explained Sample Var:

Suppose the observations are standardized.

Canonical coefficients are  $\hat{A}_z, \hat{B}_z$ .

$$\Rightarrow \begin{cases} \text{Cov}(Z^{(1)}, \hat{u}_z) = (\hat{A}_z^T)^{-1} = e^2(Z^{(1)}, \hat{u}_z) \\ \text{Cov}(Z^{(2)}, \hat{v}_z) = (\hat{B}_z^T)^{-1} = e^2(Z^{(2)}, \hat{v}_z) \end{cases}$$

$$\Rightarrow \begin{cases} \text{tr}(R_{11}) = \text{tr}\left(\sum_{i=1}^p \hat{a}_z^{(i)} \hat{a}_z^{(i)T}\right) = p = \sum_{i=1}^r \sum_{j=1}^p \hat{r}_{\hat{u}_{z1}, \hat{z}_k^{(1)}}^2 \\ \text{tr}(R_{22}) = \text{tr}\left(\sum_{i=1}^p \hat{b}_z^{(i)} \hat{b}_z^{(i)T}\right) = q = \sum_{i=1}^r \sum_{j=1}^p \hat{r}_{\hat{v}_{z1}, \hat{z}_k^{(2)}}^2 \end{cases}$$

Impk: We can calculate the proportion of total sample variances explained by first  $r$  variables.

$$R_{Z^{(1)} | \hat{u}_{z1}, \dots, \hat{u}_{zr}}^2 = \text{tr}\left(\sum_{i=1}^r \hat{a}_z^{(i)} \hat{a}_z^{(i)T}\right) / \text{tr}(R_{11}) = \frac{\sum_{i=1}^r \sum_{j=1}^p \hat{r}_{\hat{u}_{z1}, \hat{z}_k^{(1)}}^2}{p}$$

It indicates how well the  $r$  canonical variates represent the original sets  $Z^{(1)}$ .

And the matrix of error can be interpreted by

$1 - R_{Z^{(1)} | \hat{u}_{z1}, \dots, \hat{u}_{zr}}^2$ . Similar for  $Z^{(2)}, \hat{v}_z$ .

## (4) Large Sample Test:



Assume:  $\begin{pmatrix} X_j^{(1)} \\ X_j^{(2)} \end{pmatrix} \stackrel{i.i.d}{\sim} N_{p+q}(\mu, \Sigma), 1 \leq j \leq n.$

Test  $H_0: \Sigma_{12} = 0$  v.s.  $H_1: \Sigma_{12} \neq 0$ . By MLE test:

$$-2 \ln \Lambda = n \ln \left( \frac{|S_{11}| |S_{22}|}{|S|} \right) = n \ln |I - S_{11}^{-1} S_{12} S_{22}^{-1} S_{21}|$$

$$\begin{aligned} \text{By } \begin{vmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{vmatrix} &= \begin{vmatrix} I & 0 \\ -S_{21} S_{11}^{-1} & I \end{vmatrix} \begin{vmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{vmatrix} \begin{vmatrix} I & -S_{11}^{-1} S_{12} \\ 0 & I \end{vmatrix} \\ &= |S_{22}| |S_{11} - S_{12} S_{22}^{-1} S_{21}| \stackrel{A}{=} |S_{22}| |S_{11}| |I - \hat{M}| \end{aligned}$$

$$\therefore \Lambda = \frac{1}{n} \prod_{i=1}^p (1 - \hat{e}_i^{*2}). \quad T = -n \sum_{i=1}^p \ln(1 - \hat{e}_i^{*2})$$

$$T \sim \chi_{p2}^2 \text{ when } n \rightarrow \infty.$$

Note:  $T \uparrow$  if  $\exists \hat{e}_i^{*2} \rightarrow 1$ .  $\therefore R = \{T > \chi_{p2}^2(\alpha)\}$ .

Remark: i)  $H_0: \Sigma_{12} = 0 \Leftrightarrow H_0: e_1^* = \dots = e_p^* = 0$

ii) Bartlett suggest:  $-(n-1 - \frac{1}{2}(p+q+1)) \ln \Lambda$

Continuation: Test:  $H_0^k: e_1^* \dots e_k^* \neq 0, e_{k+1}^* \dots e_p^* = 0$

v.s.  $H_1^k: \exists i \geq k+1, e_i^* \neq 0$

$$R = \left\{ -(n-1 - \frac{1}{2}(p+q+1)) \ln \frac{1}{n} \prod_{i=k+1}^p (1 - \hat{e}_i^{*2}) > \chi_{(q-k)(q-k)}^2(\alpha) \right\}$$

Remark: It can reduce the numbers of canonical variables. (c.f. (3), (4))

(5) Procedure:

i) From samples of  $X, Y \Rightarrow$  Calculate  $R$

ii) Calculate canonical coefficients and variables

iii) Test  $\Sigma_{xy} = 0$ ?

iv) Apply CCA based on ii)