

Multivariate Analysis

Canonical Correlation Analysis

Shaobo Jin

Department of Mathematics

Intended Learning Outcome

Through this chapter, you should be able to

- ① perform canonical correlation analysis,
- ② perform partial least squares regression.

Motivation

- PCA performs dimension reduction for one set of p variables, while extracting as much information as possible.
- Now suppose that we have a $p \times 1$ vector of variables $\mathbf{X}^{(1)}$, and a $q \times 1$ vector of variables $\mathbf{X}^{(2)}$, measured on the same individual. We want to summarize the relationships between two vectors. '
 - e.g., you have p variables ($\mathbf{X}^{(1)}$) to measure an individual's motivations for watching online videos, and q variables ($\mathbf{X}^{(2)}$) to measure how the individual access the online videos.
 - pq pairwise scatter plots.
 - **Canonical correlation analysis (CCA)** focuses on linear combinations of variables such that much fewer plots are needed.

Task

Let $\mathbf{X}^{(1)}$ be a $p \times 1$ random vector, and $\mathbf{X}^{(2)}$ be a $q \times 1$ random vector. We assume $p \leq q$ with loss of generality. Let

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}^{(1)} \\ \mathbf{X}^{(2)} \end{bmatrix},$$

where

$$\mathbb{E}(\mathbf{X}) = \begin{bmatrix} \boldsymbol{\mu}^{(1)} \\ \boldsymbol{\mu}^{(2)} \end{bmatrix}, \quad \text{cov}(\mathbf{X}) = \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{bmatrix}.$$

The main task of CCA is to find linear combinations

$$U = \mathbf{a}^T \mathbf{X}^{(1)}, \quad V = \mathbf{b}^T \mathbf{X}^{(2)},$$

such that $\text{corr}(U, V)$ is maximized.

Correlation Coefficient

By [Result 4.2](#),

$$\begin{aligned}
 \text{cov} \left(\begin{bmatrix} U \\ V \end{bmatrix} \right) &= \text{cov} \left(\begin{bmatrix} \mathbf{a}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{b}^T \end{bmatrix} \begin{bmatrix} \mathbf{X}^{(1)} \\ \mathbf{X}^{(2)} \end{bmatrix} \right) \\
 &= \begin{bmatrix} \mathbf{a}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{b}^T \end{bmatrix} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} \begin{bmatrix} \mathbf{a} & \mathbf{0} \\ \mathbf{0} & \mathbf{b} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{a}^T \Sigma_{11} \mathbf{a} & \mathbf{a}^T \Sigma_{12} \mathbf{b} \\ \mathbf{b}^T \Sigma_{12}^T \mathbf{a} & \mathbf{b}^T \Sigma_{22} \mathbf{b} \end{bmatrix}.
 \end{aligned}$$

We shall seek coefficient vectors \mathbf{a} and \mathbf{b} such that

$$\text{corr}(U, V) = \frac{\mathbf{a}^T \Sigma_{12} \mathbf{b}}{\sqrt{\mathbf{a}^T \Sigma_{11} \mathbf{a}} \sqrt{\mathbf{b}^T \Sigma_{22} \mathbf{b}}}$$

is as large as possible.

Restriction

Note that

$$\begin{aligned}\text{corr}(U, V) &= \frac{\mathbf{a}^T \boldsymbol{\Sigma}_{12} \mathbf{b}}{\sqrt{\mathbf{a}^T \boldsymbol{\Sigma}_{11} \mathbf{a}} \sqrt{\mathbf{b}^T \boldsymbol{\Sigma}_{22} \mathbf{b}}} \\ s &= \frac{(\mathbf{c}\mathbf{a})^T \boldsymbol{\Sigma}_{12} (\mathbf{c}\mathbf{b})}{\sqrt{(\mathbf{c}\mathbf{a})^T \boldsymbol{\Sigma}_{11} (\mathbf{c}\mathbf{a})} \sqrt{(\mathbf{c}\mathbf{b})^T \boldsymbol{\Sigma}_{22} (\mathbf{c}\mathbf{b})}}.\end{aligned}$$

Hence, we need to set the scales of \mathbf{a} and \mathbf{b} . One option is

$$\begin{aligned}\mathbf{a}^T \boldsymbol{\Sigma}_{11} \mathbf{a} &= 1, \\ \mathbf{b}^T \boldsymbol{\Sigma}_{22} \mathbf{b} &= 1.\end{aligned}$$

Canonical Variates

- 1 The **first canonical variate pair** is the pair of linear combinations U_1 and V_1 having unit variances which maximizes the correlation $\text{corr}(U_1, V_1)$ (**first canonical correlation**).
- 2 The **second canonical variate pair** is the pair of linear combinations U_2 and V_2 having unit variances which maximizes the correlation $\text{corr}(U_2, V_2)$ (**second canonical correlation**) among all choices that are uncorrelated with the first pair of canonical variables.
- 3 The **k th canonical variate pair** is the pair of linear combinations U_k and V_k having unit variances which maximizes the correlation $\text{corr}(U_k, V_k)$ (**k th canonical correlation**) among all choices that are uncorrelated with the previous $k - 1$ pairs of canonical variables.

A Useful Lemma

Lemma

If λ is an eigenvalue of \mathbf{AB} with corresponding eigenvector \mathbf{x} , i.e.,

$$\mathbf{ABx} = \lambda \mathbf{x}.$$

Then, λ is an eigenvalue of \mathbf{BA} with corresponding eigenvector \mathbf{Bx} , since

$$\mathbf{BABx} = \lambda \mathbf{Bx}.$$

Find Canonical Variates

Result 10.1

Suppose $p \leq q$ and let the random vectors $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ have a full rank covariance matrix Σ . Then the k th canonical correlation is ρ_k^* , and the k th canonical variate pair is attained by

$$\begin{aligned}U_k &= \left(\Sigma_{11}^{-1/2} \mathbf{e}_k \right)^T \mathbf{X}^{(1)}, \\V_k &= \left(\Sigma_{22}^{-1/2} \mathbf{f}_k \right)^T \mathbf{X}^{(2)}.\end{aligned}$$

Here $\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$ are the eigenvalues of $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}^T \Sigma_{11}^{-1/2}$, and \mathbf{e}_k are the associated eigenvectors. Each \mathbf{f}_k is proportional to $\Sigma_{22}^{-1/2} \Sigma_{12}^T \Sigma_{11}^{-1/2} \mathbf{e}_k$.

Unit Variance and Zero Correlation

The canonical variates have the properties

$$\begin{aligned}\text{var}(U_k) &= \text{var}(V_k) = 1, \\ \text{cov}(U_k, U_\ell) &= 0, \quad k \neq \ell, \\ \text{cov}(V_k, V_\ell) &= 0, \quad k \neq \ell, \\ \text{cov}(U_k, V_\ell) &= 0, \quad k \neq \ell.\end{aligned}$$

for all $k, \ell = 1, 2, \dots, p$.

If we apply PCA to $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ separately, then we only have

$$\begin{aligned}\text{cov}(U_k, U_\ell) &= 0, \quad k \neq \ell, \\ \text{cov}(V_k, V_\ell) &= 0, \quad k \neq \ell.\end{aligned}$$

Scale Invariant: Canonical Correlation

Suppose that we change the scale as $\mathbf{Z}^{(1)} = \mathbf{C}_1 \mathbf{X}^{(1)} + \mathbf{d}_1$ and $\mathbf{Z}^{(2)} = \mathbf{C}_2 \mathbf{X}^{(2)} + \mathbf{d}_2$, where \mathbf{C}_1 and \mathbf{C}_2 are invertible. Then,

$$\text{cov} \left(\begin{bmatrix} \mathbf{Z}^{(1)} \\ \mathbf{Z}^{(2)} \end{bmatrix} \right) = \begin{bmatrix} \mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T & \mathbf{C}_1 \boldsymbol{\Sigma}_{12} \mathbf{C}_2^T \\ \mathbf{C}_2 \boldsymbol{\Sigma}_{12}^T \mathbf{C}_1^T & \mathbf{C}_2 \boldsymbol{\Sigma}_{22} \mathbf{C}_2^T \end{bmatrix}.$$

The k th canonical correlation is r_k^* , the eigenvalues of

$$\begin{aligned} & (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{C}_1 \boldsymbol{\Sigma}_{12} \mathbf{C}_2^T (\mathbf{C}_2 \boldsymbol{\Sigma}_{22} \mathbf{C}_2^T)^{-1} (\mathbf{C}_1 \boldsymbol{\Sigma}_{12} \mathbf{C}_2^T)^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \\ &= (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{C}_1 \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2}. \end{aligned}$$

They are the same as the nonzero eigenvalues of

$$\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \boldsymbol{\Sigma}_{11}^{-1/2}.$$

Scale Invariant: Coefficient Vector

Let $\mathbf{e}_k^{(Z)}$ satisfies

$$(\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{C}_1 \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} = \rho_k^* \mathbf{e}_k^{(Z)}.$$

Then,

$$\begin{aligned} (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1} \mathbf{C}_1 \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)} &= \rho_k^* (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)} \\ \Rightarrow (\mathbf{C}_1^T)^{-1} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)} &= \rho_k^* (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)} \\ \Rightarrow \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T \times \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)} &= \rho_k^* \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)}. \end{aligned}$$

This means that

$$\boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_k = \mathbf{C}_1^T (\mathbf{C}_1 \boldsymbol{\Sigma}_{11} \mathbf{C}_1^T)^{-1/2} \mathbf{e}_k^{(Z)}.$$

Canonical Correlation As Correlation Bound

Note that

$$\left| \text{cor} \left(X_i^{(1)}, X_j^{(j)} \right) \right| \leq \max_{a,b} \left| \text{corr} \left(\mathbf{a}^T \mathbf{X}^{(1)}, \mathbf{b}^T \mathbf{X}^{(2)} \right) \right| = \rho_1^*.$$

Hence, the first canonical correlation is no lower than the absolute value of the correlation between any $X_i^{(1)}$ and $X_j^{(j)}$.

Proportion of Explained Variance

In fact,

$$\begin{aligned}\max_b \operatorname{corr} \left(U_k, \mathbf{b}^T \mathbf{X}^{(2)} \right) &= \rho_k^*, \\ \max_a \operatorname{corr} \left(\mathbf{a}^T \mathbf{X}^{(1)}, V_k \right) &= \rho_k^*,\end{aligned}$$

for any k .

The k th squared canonical correlation ρ_k^{*2} is

- the proportion of the variance of canonical variate U_k explained by the set $\mathbf{X}^{(2)}$,
- the proportion of the variance of canonical variate V_k explained by the set $\mathbf{X}^{(1)}$.

Hence, ρ_k^{*2} is the shared variance between $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$.

Sample Canonical Variate Pair

A random sample of n observations is assembled into the $n \times (p + q)$ data matrix $\mathbf{X} = [\mathbf{X}^{(1)} \quad \mathbf{X}^{(2)}]$.

- 1 The **first sample canonical variate pair** is the pair of linear combinations \hat{U}_1 and \hat{V}_1 having unit variances which maximizes the sample correlation $r_{\hat{U}_1, \hat{V}_1}$.
- 2 The **second sample canonical variate pair** is the pair of linear combinations \hat{U}_2 and \hat{V}_2 having unit variances which maximizes the sample correlation $r_{\hat{U}_2, \hat{V}_2}$ among all choices that are uncorrelated with the first sample canonical variate pair.
- 3 The **k th sample canonical variate pair** is the pair of linear combinations \hat{U}_k and \hat{V}_k having unit variances which maximizes the sample correlation $r_{\hat{U}_k, \hat{V}_k}$ among all choices that are uncorrelated with the previous $k - 1$ sample canonical variate pairs.

Sample Correlation

Let

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{12}^T & \mathbf{S}_{22} \end{bmatrix}$$

be the sample covariance matrix of \mathbf{X} . The sample correlation between $\mathbf{a}^T \mathbf{X}^{(1)}$ and $\mathbf{b}^T \mathbf{X}^{(2)}$ is

$$\frac{\mathbf{a}^T \mathbf{S}_{12} \mathbf{b}}{\sqrt{\mathbf{a}^T \mathbf{S}_{11} \mathbf{a}} \sqrt{\mathbf{b}^T \mathbf{S}_{22} \mathbf{b}}}.$$

The population correlation is

$$\frac{\mathbf{a}^T \boldsymbol{\Sigma}_{12} \mathbf{b}}{\sqrt{\mathbf{a}^T \boldsymbol{\Sigma}_{11} \mathbf{a}} \sqrt{\mathbf{b}^T \boldsymbol{\Sigma}_{22} \mathbf{b}}}.$$

Find Sample Canonical Variate Pair

Result 10.2

Let $\hat{\rho}_1^{*2} \geq \hat{\rho}_2^{*2} \geq \dots \geq \hat{\rho}_p^{*2}$ be the eigenvalues of $\mathbf{S}_{11}^{-1/2} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1/2}$ with corresponding eigenvectors $\hat{\mathbf{e}}_1, \dots, \hat{\mathbf{e}}_p$, where $p \leq q$. Let $\hat{\mathbf{f}}_1, \dots, \hat{\mathbf{f}}_p$ be the eigenvectors of $\mathbf{S}_{22}^{-1/2} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-1/2}$ where the first p eigenvectors may be obtained from $\hat{\mathbf{f}}_k = \hat{\rho}_k^{*-1} \mathbf{S}_{22}^{-1/2} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1} \hat{\mathbf{e}}_k$, $k = 1, 2, \dots, p$. Then the k th sample canonical variate pair is

$$\hat{U}_k = \left(\mathbf{S}_{11}^{-1/2} \hat{\mathbf{e}}_k \right)^T \mathbf{x}^{(1)}, \quad \hat{V}_k = \left(\mathbf{S}_{22}^{-1/2} \hat{\mathbf{f}}_k \right)^T \mathbf{x}^{(2)}.$$

The k th sample canonical correlation is $\hat{\rho}_k^*$.

Unit Variance and Zero Correlation

The sample canonical variates have the properties

$$\begin{aligned}S_{\hat{U}_k} &= S_{\hat{V}_k} = 1, \\r_{\hat{U}_k, \hat{U}_\ell} &= 0, \quad k \neq \ell, \\r_{\hat{V}_k, \hat{V}_\ell} &= 0, \quad k \neq \ell, \\r_{\hat{U}_k, \hat{V}_\ell} &= 0, \quad k \neq \ell,\end{aligned}$$

for all $k, \ell = 1, 2, \dots, p$.

If we apply PCA to $\mathbf{X}^{(1)}$ and $\mathbf{X}^{(2)}$ separately, then we only have

$$\begin{aligned}r_{\hat{U}_k, \hat{U}_\ell} &= 0, \quad k \neq \ell, \\r_{\hat{V}_k, \hat{V}_\ell} &= 0, \quad k \neq \ell.\end{aligned}$$

Special Case: $p = 1$

Consider a special case where $p = 1$. We denote $Y = X^{(1)}$ and $\mathbf{Z} = \mathbf{X}^{(2)}$. The matrix $\mathbf{S}_{11}^{-1/2} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1/2}$ reduces to a scalar with eigenvalue itself and eigenvector 1. Then, the canonical variate pair is attained by

$$\begin{aligned} V &= \left(\mathbf{S}_{22}^{-1/2} \mathbf{f} \right)^T \mathbf{Z} = \left(\mathbf{S}_{22}^{-1/2} \times |\rho^*| \mathbf{S}_{22}^{-1/2} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1/2} \hat{e} \right)^T \mathbf{Z} \\ &= \mathbf{S}_{11}^{-1} \sqrt{\mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{12}^T \mathbf{S}_{12} \mathbf{S}_{22}^{-1}} \mathbf{Z}, \end{aligned}$$

where $\hat{\rho}^{*2} = \mathbf{S}_{11}^{-1/2} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{12}^T \mathbf{S}_{11}^{-1/2} = \mathbf{S}_{11}^{-1} \mathbf{S}_{12} \mathbf{S}_{22}^{-1} \mathbf{S}_{12}^T$ and $\hat{e} = 1$.

If we apply classic linear regression (regress Y on \mathbf{Z}), the OLS fitted model is

$$\hat{Y} = \underbrace{\mathbf{y}^T \mathbf{Z}_D}_{\propto \mathbf{S}_{12}} \underbrace{(\mathbf{Z}_D^T \mathbf{Z}_D)^{-1}}_{\propto \mathbf{S}_{22}^{-1}} \mathbf{Z},$$

where \mathbf{Z}_D is the demeaned data matrix of \mathbf{Z} and \mathbf{y} is the demeaned response vector of Y .

CCA Versus Regression: Seemingly Different

```
## Canonical Correlation Analysis
CC <- cc(X = as.matrix(Data[, "Y"]),
        Y = as.matrix(Data[, c("Z1", "Z2", "Z3")]))
CC$ycoef

##           [,1]
## Z1 0.6444941
## Z2 0.3009713
## Z3 0.3615434

## Classic Linear Regression
LM <- lm(Y ~ Z1 + Z2 + Z3, data = Data)
coef(LM)

## (Intercept)           Z1           Z2           Z3
## 0.1486136    0.2558822    0.1194941    0.1435429
```

More Outputs of Linear Regression

```
summary(LM)

##
## Call:
## lm(formula = Y ~ Z1 + Z2 + Z3, data = Data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.30639 -0.61481  0.05427  0.56709  1.90861
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.14861    0.09000   1.651   0.1019
## Z1           0.25588    0.09957   2.570   0.0117 *
## Z2           0.11949    0.09509   1.257   0.2119
## Z3           0.14354    0.09993   1.436   0.1541
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8756 on 96 degrees of freedom
## Multiple R-squared:  0.1749, Adjusted R-squared:  0.1492
## F-statistic: 6.785 on 3 and 96 DF,  p-value: 0.0003378
```

Almost Same Thing

```
## Just Scaling  
c(CC$ycoef / coef(LM)[-1])  
  
## [1] 2.518714 2.518714 2.518714  
  
## Canonical Correlation is related  
c(CC$cor ^ 2, summary(LM)$r.squared)  
  
## [1] 0.1749436 0.1749436
```

Alternative to Multivariate Multiple Regression

Suppose that we have p response variables and q covariates that we can use to model the response variable.

- 1 One approach is to use multivariate multiple regression (Chapter 7).
- 2 An alternative approach is **partial least squares** regression (**PLS** regression).

A little history:

- Karl Gustav Jöreskog popularized factor analysis. He had professorship at Uppsala 1971–2000.
- Herman Wold developed PLS regression. He had professorship at Uppsala 1942–1970.
- Jöreskog was a student of Wold at Uppsala.

Maximize Covariance

PLS regression is an alternative to multivariate multiple regression if

- n is relatively small comparing to p or q ,
- correlation between covariates are high (multicollinearity).

Suppose that we have demeaned the p response variables and q covariates, such that the sample mean of each variable is 0. We want to find linear combinations $T = \mathbf{a}^T \mathbf{X}$ and $U = \mathbf{b}^T \mathbf{Y}$ such that the covariance between T_k and U is maximized. That is,

$$\begin{aligned} \max_{\mathbf{a}, \mathbf{b}} \text{cov}(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}) \\ \text{s.t. } \mathbf{a}^T \mathbf{a} = \mathbf{b}^T \mathbf{b} = 1. \end{aligned}$$

This is similar to the objective of CCA.

Eigenvalue and Eigenvector

Since the variables are demeaned, the sample covariance matrix between $\mathbf{X}_{n \times q}$ and $\mathbf{Y}_{n \times p}$ can be computed by $n^{-1}\mathbf{X}^T\mathbf{Y}$ or $(n-1)^{-1}\mathbf{X}^T\mathbf{Y}$. Then,

$$\max_{\mathbf{a}, \mathbf{b}} \mathbf{a}^T \mathbf{X}^T \mathbf{Y} \mathbf{b} \quad \text{s.t. } \mathbf{a}^T \mathbf{a} = \mathbf{b}^T \mathbf{b} = 1.$$

Hence, \mathbf{a} is the eigenvector of $\mathbf{X}^T \mathbf{Y} \mathbf{Y}^T \mathbf{X}$ corresponding to the largest eigenvalue, and $\mathbf{b} \propto \mathbf{Y}^T \mathbf{X} \mathbf{a}$.

- 1 We regress \mathbf{X} on $\mathbf{t} = \mathbf{X} \mathbf{a}$, a $n \times 1$ vector. The OLS estimator of regression coefficients is $\hat{\boldsymbol{\beta}} = (\mathbf{t}^T \mathbf{t})^{-1} \mathbf{t}^T \mathbf{X}$.
- 2 We regress \mathbf{Y} also on \mathbf{t} . The OLS estimator of regression coefficients is $\hat{\boldsymbol{\gamma}} = (\mathbf{t}^T \mathbf{t})^{-1} \mathbf{t}^T \mathbf{Y}$.

PLS Regression

From the regression model of \mathbf{Y} on \mathbf{t} , our fitted model for a new value t_0 satisfies

$$\hat{\mathbf{y}} = \hat{\gamma}t_0.$$

Since $t_0 = \mathbf{a}^T \mathbf{x}_0$, then

$$\hat{\mathbf{y}} = \hat{\gamma}\mathbf{a}^T \mathbf{x}_0,$$

which is the **PLS regression** of \mathbf{y} on \mathbf{x} .

More Components

So far we have only used one pair of linear combinations. Similar to CCA, we can extract more components.

- When we regress \mathbf{X} on $\mathbf{t} = \mathbf{X}\mathbf{a}$, the residual matrix is

$$\mathbf{E}_X = \mathbf{X} - \mathbf{t}(\mathbf{t}^T \mathbf{t})^{-1} \mathbf{t}^T \mathbf{X}.$$

- We regress \mathbf{Y} also on \mathbf{t} , the residual matrix is

$$\mathbf{E}_Y = \mathbf{Y} - \mathbf{t}(\mathbf{t}^T \mathbf{t})^{-1} \mathbf{t}^T \mathbf{Y}.$$

- We treat \mathbf{E}_X and \mathbf{E}_Y as \mathbf{X} and \mathbf{Y} respectively in a new iteration of PLS regression.