

Multivariate Regression

This note is prepared based on *Chapter 6, Multivariate Regression* in Izenman (2009).

I. Introduction

1. **Overview:** *Multivariate regression* is an extension of the multiple regression and has s output variables $Y = (Y_1, \dots, Y_s)^\top \in \mathbb{R}^s$, each of whose behavior may be influenced by exactly the same set of predictors $X = (X_1, \dots, X_p)^\top \in \mathbb{R}^p$.
2. **Goal:** We are interested in estimating the regression relationship between Y and X , taking into account the various *dependencies* between the p -dimensional vector X and the s -dimensional vector Y and the dependencies within X and within Y .

Remark. In the multivariate regression setting,

- (a) the components of X are correlated with each other,
- (b) the components of Y are correlated with each other, and
- (c) the components of Y are correlated with components of X .

II. Random Design Case

1. **Setup:** Assume that the random vectors

$$X = (X_1, \dots, X_p)^\top \quad \text{and} \quad Y = (Y_1, \dots, Y_s)^\top$$

are jointly distributed, where the mean vectors of X and Y are μ_X and μ_Y , respectively, and the joint covariance matrix

$$\begin{pmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{pmatrix}.$$

In addition, we assume that $s \leq p$.

2. **Classical Multivariate Regression Model:** Assume that Y is related to X by the following *multivariate linear model*

$$Y = \mu + \Theta X + \varepsilon, \tag{1}$$

where $\mu \in \mathbb{R}^s$ and $\Theta \in \mathbb{R}^{s \times p}$ are the unknown parameters to be estimated, and $\varepsilon \in \mathbb{R}^s$ is the unobservable error component of the model with mean $\mathbf{0}_s$ and unknown error covariance matrix $\text{Cov}[\varepsilon] = \Sigma_{\varepsilon\varepsilon}$. We assume that X and ε are independent.

3. Least Squares Estimate of Parameters: We are interested in finding μ and Θ that minimize

$$L(\mu, \Theta) := \mathbb{E}[(Y - \mu - \Theta X)(Y - \mu - \Theta X)^\top], \quad (2)$$

where the expectation is taken over the joint distribution of X and Y .

(a) *First Derivation:* Let $Y_c := Y - \mu_Y$ and $X_c := X - \mu_X$ and assume $\Sigma_{XX} \succ 0$. We have

$$\begin{aligned} L(\mu, \Theta) &= \mathbb{E}[(Y - \mu - \Theta X)(Y - \mu - \Theta X)^\top] \\ &= \mathbb{E}[Y_c Y_c^\top - Y_c X_c^\top \Theta^\top - \Theta X_c Y_c^\top + \Theta X_c X_c^\top \Theta^\top] \\ &\quad + (\mu - \mu_Y + \Theta \mu_X)(\mu - \mu_Y + \Theta \mu_X)^\top \\ &= \Sigma_{YY} - \Sigma_{YX} \Theta^\top - \Theta \Sigma_{XY} + \Theta \Sigma_{XX} \Theta^\top \\ &\quad + (\mu - \mu_Y + \Theta \mu_X)(\mu - \mu_Y + \Theta \mu_X)^\top. \end{aligned}$$

We note that

$$\begin{aligned} & - \Sigma_{YX} \Theta^\top - \Theta \Sigma_{XY} + \Theta \Sigma_{XX} \Theta^\top \\ &= (\Theta \Sigma_{XX}^{1/2} - \Sigma_{YX} \Sigma_{XX}^{-1/2})(\Theta \Sigma_{XX}^{1/2} - \Sigma_{YX} \Sigma_{XX}^{-1/2})^\top - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}. \end{aligned}$$

Hence,

$$\begin{aligned} L(\mu, \Theta) &= (\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}) + (\Theta \Sigma_{XX}^{1/2} - \Sigma_{YX} \Sigma_{XX}^{-1/2})(\Theta \Sigma_{XX}^{1/2} - \Sigma_{YX} \Sigma_{XX}^{-1/2})^\top \\ &\quad + (\mu - \mu_Y + \Theta \mu_X)(\mu - \mu_Y + \Theta \mu_X)^\top \\ &\geq \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}, \end{aligned}$$

with equality being held when

$$\Theta^* := \Sigma_{YX} \Sigma_{XX}^{-1}, \quad \mu^* := \mu_Y - \Theta \mu_X.$$

In other words,

$$(\mu^*, \Theta^*) = \arg \min_{\mu, \Theta} L(\mu, \Theta).$$

(b) *Second Derivation:* Assuming that we can interchange the derivative and the integral, we take the partial derivatives of $L(\mu, \Theta)$ with respect to the two arguments and set the derivatives to zero:

$$\begin{aligned} \frac{\partial L(\mu, \Theta)}{\partial \mu} &= -2\mu_Y^\top + 2\mu^\top + 2\mu_X \Theta^\top \stackrel{\text{set}}{=} \mathbf{0}_s, \\ \frac{\partial L(\mu, \Theta)}{\partial \Theta} &= -2 \mathbb{E}[(Y - \mu - \Theta X)X^\top] \stackrel{\text{set}}{=} \mathbf{0}_{s \times p}. \end{aligned}$$

It follows that $(\mu^*, \Theta^*) = \arg \min_{\mu, \Theta} L(\mu, \Theta)$ must satisfy

$$\mu^* = \mu_Y - \Theta^* \mu_X, \quad \text{and} \quad \Theta^* = \Sigma_{YX} \Sigma_{XX}^{-1}.$$

Here, Θ^* is the *full-rank regression coefficient matrix* of Y on X , and

$$Y = \mu_Y + \Sigma_{YX} \Sigma_{XX}^{-1} (X - \mu_X)$$

is the *full-rank linear regression function* of Y on X . Here, the “full-rank” refers to the rank of Θ . At the minimum of L , the error is

$$\varepsilon = Y - \mu_Y - \Sigma_{YX} \Sigma_{XX}^{-1} (X - \mu_X) = Y_c - \Sigma_{YX} \Sigma_{XX}^{-1} X_c.$$

In particular, note that

$$\begin{aligned} \mathbb{E}[\varepsilon] &= \mathbf{0}_s, \\ \text{Var}[\varepsilon] &= \Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}, \\ \mathbb{E}[\varepsilon X_c^\top] &= \mathbf{0}_{s \times p}. \end{aligned}$$

4. Multivariate Reduced-Rank Regression:

(a) *Model Specification:* Consider the multivariate linear regression model given by

$$Y = \mu + CX + \varepsilon, \quad (3)$$

where $\mu \in \mathbb{R}^s$ and $C \in \mathbb{R}^{s \times p}$ are unknown regression parameters and $\varepsilon \in \mathbb{R}^s$ is the unobservable error with mean $\mathbb{E}[\varepsilon] = \mathbf{0}_s$ and covariance matrix $\text{Cov}[\varepsilon] = \Sigma_{\varepsilon\varepsilon}$. We assume that ε and X are independent. Here, we allow the rank of the regression coefficient matrix C to be deficient, i.e.,

$$\text{rank}(C) = t \leq \min s = \{s, p\}. \quad (4)$$

Remark 1. The “reduced-rank” condition (4) on C brings a true multivariate feature into the model, implying that there may be a number of *linear constraints* on the set of regression coefficients in the model.

Remark 2. The name *reduced-rank regression* is used to distinguish the case $1 \leq t < s$ from the *full-rank regression* with $t = s$.

(b) *Question of Interest:* When $\text{rank}(C) = t$, there exists two full-rank matrices, an $s \times t$ matrix A and a $t \times p$ matrix B , such that $C = AB$. Note that this decomposition is *not* unique since we can always find a nonsingular $t \times t$ matrix T such that

$$C = AB = (AT)(T^{-1}B) = DE.$$

With the decomposition $C = AB$, we can write (3) as

$$Y = \mu + ABX + \varepsilon.$$

We wish to estimate the unknown parameters A , B and μ .

(c) *Effective Dimensionality:* The rank of the matrix C here, t , is a meta-parameter called the *effective dimensionality* of the multivariate regression.

- (d) *Review of the Eckart-Young's Inequality:* Suppose both \mathbf{A} and \mathbf{B} are matrices of size $m \times n$. Assume that \mathbf{B} has the reduced rank $\text{rank}(\mathbf{B}) = t$ and \mathbf{A} is of full rank, $\text{rank}(\mathbf{A}) = \min\{m, n\}$. We use \mathbf{B} to approximate \mathbf{A} . Then,

$$\lambda_j((\mathbf{A} - \mathbf{B})(\mathbf{A} - \mathbf{B})^\top) \geq \lambda_{j+t}(\mathbf{A}\mathbf{A}^\top),$$

with equality if

$$\mathbf{B} = \sum_{i=1}^t \lambda_i^{1/2} \mathbf{u}_i \mathbf{v}_i^\top,$$

where $\lambda_i = \lambda_i(\mathbf{A}\mathbf{A}^\top)$, \mathbf{u}_i is the i -th eigenvector of $\mathbf{A}\mathbf{A}^\top$, and \mathbf{v}_i is the i -th eigenvector of $\mathbf{A}^\top \mathbf{A}$.

- (e) *Minimizing a Weighted Least Squares Criterion:* We minimize the following weighted sum-of-squares criterion

$$L(t) := \mathbb{E}[(Y - \boldsymbol{\mu} - \mathbf{A}\mathbf{B}\mathbf{X})^\top \boldsymbol{\Gamma} (Y - \boldsymbol{\mu} - \mathbf{A}\mathbf{B}\mathbf{X})], \quad (5)$$

where $\boldsymbol{\Gamma} \in \mathbb{R}^{s \times s}$ is a positive definite symmetric matrix of weights and the expectation is taken over the joint distribution of (X, Y) . By letting $Y_c := Y - \boldsymbol{\mu}_Y$ and $X_c := X - \boldsymbol{\mu}_X$, first note that

$$\begin{aligned} L(t) &\geq \mathbb{E}[(Y_c - \mathbf{C}X_c)^\top \boldsymbol{\Gamma} (Y_c - \mathbf{C}X_c)] \\ &= \text{trace}(\boldsymbol{\Sigma}_{YY}^* - \mathbf{C}^* \boldsymbol{\Sigma}_{XY}^* - \boldsymbol{\Sigma}_{YX}^* \mathbf{C}^* + \mathbf{C}^* \boldsymbol{\Sigma}_{XX}^* \mathbf{C}^*) \\ &= \text{trace}((\boldsymbol{\Sigma}_{YY}^* - \boldsymbol{\Sigma}_{YX}^* \boldsymbol{\Sigma}_{XX}^{*-1} \boldsymbol{\Sigma}_{XY}^*)) \\ &\quad + \text{trace}((\mathbf{C}^* \boldsymbol{\Sigma}_{XX}^{*1/2} - \boldsymbol{\Sigma}_{YX}^* \boldsymbol{\Sigma}_{XX}^{*-1/2})(\mathbf{C}^* \boldsymbol{\Sigma}_{XX}^{*1/2} - \boldsymbol{\Sigma}_{YX}^* \boldsymbol{\Sigma}_{XX}^{*-1/2})^\top), \end{aligned} \quad (6)$$

where the last equality follows by completing the perfect square, and $\boldsymbol{\Sigma}_{XX}^* := \boldsymbol{\Sigma}_{XX}$, $\boldsymbol{\Sigma}_{YY}^* := \boldsymbol{\Gamma}^{1/2} \boldsymbol{\Sigma}_{YY} \boldsymbol{\Gamma}^{1/2}$, $\boldsymbol{\Sigma}_{XY}^* := \boldsymbol{\Sigma}_{XY} \boldsymbol{\Gamma}^{1/2}$, and $\mathbf{C}^* := \boldsymbol{\Gamma}^{1/2} \mathbf{C}$.

Now, we assume that $\text{rank}(\mathbf{C}) = t$. According to the Eckart-Young's Inequality, the second trace is minimized when

$$\mathbf{C}^* \boldsymbol{\Sigma}_{XX}^{*1/2} = \sum_{i=1}^t \lambda_i^{1/2} \mathbf{v}_i \mathbf{u}_i^\top,$$

where λ_i is the i -th eigenvalue of $\boldsymbol{\Sigma}_{YX}^* \boldsymbol{\Sigma}_{XX}^{*-1} \boldsymbol{\Sigma}_{XY}^* = \boldsymbol{\Gamma}^{1/2} \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \boldsymbol{\Gamma}^{1/2}$, and \mathbf{v}_i is the eigenvector associated with λ_i , and

$$\mathbf{u}_i = \lambda_i^{-1/2} \boldsymbol{\Sigma}_{XX}^{*-1/2} \boldsymbol{\Sigma}_{XY}^* \mathbf{v}_i = \lambda_i^{-1/2} \boldsymbol{\Sigma}_{XX}^{-1/2} \boldsymbol{\Sigma}_{XY} \boldsymbol{\Gamma}^{1/2} \mathbf{v}_i.$$

It follows that the optimal \mathbf{C} of reduced rank t that minimizes L is

$$\mathbf{C}^{(t)} = \boldsymbol{\Gamma}^{-1/2} \left(\sum_{j=1}^t \mathbf{v}_j \mathbf{v}_j^\top \right) \boldsymbol{\Gamma}^{1/2} \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1}. \quad (7)$$

This $\mathbf{C}^{(t)}$ is called the *reduced-rank regression coefficient matrix* with rank t and weight matrix $\mathbf{\Gamma}$. It follows that L is minimized by letting $\boldsymbol{\mu}$, \mathbf{A} and \mathbf{B} be the following functions of t

$$\boldsymbol{\mu}^{(t)} = \boldsymbol{\mu}_Y - \mathbf{A}^{(t)}\mathbf{B}^{(t)}\boldsymbol{\mu}_X, \quad (8)$$

$$\mathbf{A}^{(t)} = \mathbf{\Gamma}^{-1/2}\mathbf{V}_t, \quad (9)$$

$$\mathbf{B}^{(t)} = \mathbf{V}_t^\top \mathbf{\Gamma}^{1/2} \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1}, \quad (10)$$

where $\mathbf{V}_t = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_t)$ is an $s \times t$ matrix and \mathbf{v}_j is the eigenvector associated with the j -th largest eigenvalue of the matrix $\mathbf{\Gamma}^{1/2} \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \mathbf{\Gamma}^{1/2}$.

(f) *Minimum of $W(t)$:* With $\mathbf{C}^{(t)}$ defined in (7), the minimum value of L is

$$\begin{aligned} L_{\min}(t) &= \mathbb{E}[(Y - \boldsymbol{\mu} - \mathbf{C}^{(t)}X)^\top \mathbf{\Gamma} (Y - \boldsymbol{\mu} - \mathbf{C}^{(t)}X)] \\ &= \text{trace}(\boldsymbol{\Sigma}_{YY}\mathbf{\Gamma}) - \sum_{j=1}^t \lambda_j. \end{aligned}$$

(g) *Two Remarks:*

- i. If we choose $t = s$, $\mathbf{C}^{(s)}$ reduces to the full-rank regression coefficient matrix $\boldsymbol{\Theta} = \mathbf{C}^{(s)}$;
- ii. For any t and any positive-definite matrix $\mathbf{\Gamma}$, $\mathbf{C}^{(t)}$ and $\boldsymbol{\Theta}$ are related by

$$\mathbf{C}^{(t)} = \mathbf{P}_{\mathbf{\Gamma}}^{(t)} \boldsymbol{\Theta},$$

where

$$\mathbf{P}_{\mathbf{\Gamma}}^{(t)} = \mathbf{\Gamma}^{-1/2} \left(\sum_{j=1}^t \mathbf{v}_j \mathbf{v}_j^\top \right) \mathbf{\Gamma}^{1/2}.$$

(h) *Special Cases of Reduced-Rank Regression:*

- i. If $X \equiv Y$ and $\mathbf{\Gamma} = \mathbf{I}_p$, we obtain Hotelling's principal component analysis;
 - ii. If $\mathbf{\Gamma} = \boldsymbol{\Sigma}_{YY}^{-1}$, we obtain the Hotelling's canonical correlation analysis;
 - iii. If $\mathbf{\Gamma} = \boldsymbol{\Sigma}_{YY}^{-1}$ and let Y be a vector of binary variables indicating the class belonging of observations, we obtain Fisher's linear discriminant analysis.
- (i) *Sample Estimates:* Let $\{(\mathbf{x}_i^\top, \mathbf{y}_i^\top)^\top\}_{i=1}^n$ be n i.i.d observations from $(X^\top, Y^\top)^\top$. Then,

- i. The mean vectors, $\boldsymbol{\mu}_X$ and $\boldsymbol{\mu}_Y$, can be estimated by

$$\hat{\boldsymbol{\mu}}_X = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \text{and} \quad \hat{\boldsymbol{\mu}}_Y = \bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i,$$

respectively.

ii. For all $i = 1, \dots, n$, let

$$\mathbf{x}_{c,i} = \mathbf{x}_i - \bar{\mathbf{x}}, \quad \text{and} \quad \mathbf{y}_{c,i} = \mathbf{y}_i - \bar{\mathbf{y}}$$

be the centered observations, and let

$$\mathbf{X}_c = [\mathbf{x}_{c,1}, \dots, \mathbf{x}_{c,n}] \in \mathbb{R}^{p \times n}, \quad \text{and} \quad \mathbf{Y}_c = [\mathbf{y}_{c,1}, \dots, \mathbf{y}_{c,n}] \in \mathbb{R}^{s \times n}$$

be the center data matrix.

iii. The covariance matrices, Σ_{XX} , Σ_{XY} , Σ_{YX} and Σ_{YY} , can be estimated by

$$\begin{aligned} \hat{\Sigma}_{XX} &= \frac{1}{n} \mathbf{X}_c \mathbf{X}_c^\top, \\ \hat{\Sigma}_{YX} &= \hat{\Sigma}_{XY}^\top = \frac{1}{n} \mathbf{Y}_c \mathbf{X}_c^\top, \\ \hat{\Sigma}_{YY} &= \frac{1}{n} \mathbf{Y}_c \mathbf{Y}_c^\top. \end{aligned}$$

iv. Matrices $\mathbf{A}^{(t)}$ in (9) and $\mathbf{B}^{(t)}$ in (10) can be estimated by

$$\begin{aligned} \hat{\mathbf{A}}^{(t)} &= \mathbf{\Gamma}^{-1/2} \hat{\mathbf{V}}_t, \\ \hat{\mathbf{B}}^{(t)} &= \hat{\mathbf{V}}_t^\top \mathbf{\Gamma}^{1/2} \hat{\Sigma}_{YX} \hat{\Sigma}_{XX}^{-1}, \end{aligned}$$

where

$$\hat{\mathbf{V}}_t = (\hat{\mathbf{v}}_1, \dots, \hat{\mathbf{v}}_t)$$

is an $s \times t$ -matrix with the j -th column, $\hat{\mathbf{v}}_j$, being the eigenvector associated with the j -th largest eigenvalue $\hat{\lambda}_j$ of the $s \times s$ symmetric matrix

$$\mathbf{\Gamma}^{1/2} \hat{\Sigma}_{YX} \hat{\Sigma}_{XX}^{-1} \hat{\Sigma}_{XY} \mathbf{\Gamma}^{1/2},$$

for $j = 1, 2, \dots, p$. The reduced-rank regression coefficient matrix $\mathbf{C}^{(t)}$ in (7) can be estimated by

$$\hat{\mathbf{C}}^{(t)} = \mathbf{\Gamma}^{-1/2} \left(\sum_{j=1}^t \hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^\top \right) \mathbf{\Gamma}^{1/2} \hat{\Sigma}_{YX} \hat{\Sigma}_{XX}^{-1},$$

and the full-rank regression coefficient matrix $\mathbf{\Theta}$ can be estimated by

$$\hat{\mathbf{\Theta}} = \hat{\mathbf{C}}^{(s)} = \hat{\Sigma}_{YX} \hat{\Sigma}_{XX}^{-1}.$$

v. *Estimation in the Presence of Singularity:* In the case where $\hat{\Sigma}_{XX}$ and/or $\hat{\Sigma}_{YY}$ are singular, we replace them by a slight perturbation of their diagonal entries using the idea of the ridge regression

$$\begin{aligned} \hat{\Sigma}_{XX}^{(\eta)} &= \frac{1}{n} (\mathbf{X}_c \mathbf{X}_c^\top + \eta \cdot \mathbf{I}_d), \\ \hat{\Sigma}_{YY}^{(\eta)} &= \frac{1}{n} (\mathbf{Y}_c \mathbf{Y}_c^\top + \eta \cdot \mathbf{I}_d), \end{aligned}$$

respectively, where $\eta > 0$ is a tuning parameter. Everything else, such obtaining estimates of $\mathbf{C}^{(t)}$, $\mathbf{A}^{(k)}$ and $\mathbf{B}^{(t)}$, can be proceeded as before and will depend on the choice of $\eta > 0$.

(j) *Assessing the Effective Dimensionality:* In order to choose the value of $t \in \{1, 2, \dots, s\}$, we choose the smallest integer such that the reduced-rank regression of Y on X with that integer as rank is as close (to be specified) as possible to the corresponding full-rank regression.

- i. Method 1. Let $L_{\min}(t)$ denote the minimum value of $L(t)$ for a fixed value of t . The reduction in $L_{\min}(t)$ by increasing the rank from $t = t_0$ to $t = t_1$ with $t_0 < t_1$ is

$$L_{\min}(t_0) - L_{\min}(t_1) = \sum_{j=t_0+1}^{t_1} \lambda_j,$$

which only depends on the eigenvalues of $\mathbf{\Gamma}^{1/2} \mathbf{\Sigma}_{YX} \mathbf{\Sigma}_{XX}^{-1} \mathbf{\Sigma}_{XY} \mathbf{\Gamma}^{1/2}$. Therefore, the rank of \mathbf{C} can be assessed by some monotone function of the sequence of ordered sample eigenvalues $\{\hat{\lambda}_j\}_{j=1}^s$, in which $\hat{\lambda}_j$ is compared with suitable reference values for each j , or by the sum of some monotone function of the smallest $s - t_0$ sample eigenvalues.

- ii. Method 2 — Rank Trace. Suppose the true rank of \mathbf{C} is t^* . The main idea behind *rank trace* is the following:

- A. for $1 \leq t < t^*$, the entries in both the estimated regression coefficient matrix and the residual covariance matrix change significantly each time we increase the rank;
- B. as soon as the true rank t^* is achieved, these two matrices stabilize.

The algorithm is provided in Algorithm 1.

Algorithm 1 Using Rank Trace to Assess the Effective Dimensionality of a Multivariate Regression

1: Define $\hat{\mathbf{C}}^{(0)} = \mathbf{0}_{s \times p}$ and $\hat{\Sigma}_{\epsilon\epsilon}^{(0)} = \hat{\Sigma}_{YY}$.

2: For $t = 1$ to s :

(a) compute $\hat{\mathbf{C}}^{(t)}$ and $\hat{\Sigma}_{\epsilon\epsilon}^{(t)}$, and set $\hat{\mathbf{C}}^{(s)} = \hat{\Theta}$ and $\hat{\Sigma}_{\epsilon\epsilon}^{(s)} = \hat{\Sigma}_{\epsilon\epsilon}$.

(b) compute

$$\Delta\hat{\mathbf{C}}^{(t)} = \frac{\|\hat{\Theta} - \hat{\mathbf{C}}^{(t)}\|}{\|\hat{\Theta}\|}, \quad \Delta\hat{\Sigma}_{\epsilon\epsilon}^{(t)} = \frac{\|\hat{\Sigma}_{\epsilon\epsilon} - \hat{\Sigma}_{\epsilon\epsilon}^{(t)}\|}{\|\hat{\Sigma}_{\epsilon\epsilon} - \hat{\Sigma}_{YY}\|},$$

where $\|\mathbf{A}\| = \sqrt{\text{trace}(\mathbf{A}\mathbf{A}^\top)} = (\sum_i \sum_j a_{ij}^2)^{1/2}$.

3: Make a scatterplot of the s points

$$(\Delta\hat{\mathbf{C}}^{(t)}, \Delta\hat{\Sigma}_{\epsilon\epsilon}^{(t)}),$$

for $t = 0, 1, \dots, s$, and join up successive points on the plot. This is called the *rank trace* for the multivariate reduce-rank regression of Y onto X .

4: Assess the rank of \mathbf{C} as the *smallest* rank for which both coordinates from Step 3 are approximately zero.

Note the following:

- The first point in the rank trace, corresponding to $t = 0$, is always plotted at $(1, 1)$ and the last point, corresponding to $t = s$, is always plotted at $(0, 0)$;
 - The horizontal coordinate, $\Delta\hat{\mathbf{C}}^{(t)}$, gives a quantitative representation of the difference between a reduced-rank regression coefficient matrix and its full-rank analog;
 - The vertical coordinate, $\Delta\hat{\Sigma}_{\epsilon\epsilon}^{(t)}$, shows the proportionate reduction in the residual variance matrix in using a simple full-rank model rather than the reduced-rank model.
- iii. Method 3 — Cross Validation. For each rank t , compute a sequence of estimates of prediction error using the cross validation. Identify the smallest rank such that, for larger ranks, the prediction error has stabilized and does *not* decrease significantly; this is similar to saying that at \hat{t} , there is an elbow in the plot of prediction error against the rank.

III. Fixed Design Case

1. **Assumptions:** Let $Y = (Y_1, \dots, Y_s)^\top$ be an s -dimensional random vector with mean $\boldsymbol{\mu}_Y \in \mathbb{R}^s$ and the covariance matrix $\boldsymbol{\Sigma}_{YY} \in \mathbb{R}^{s \times s}$, and let $X = (X_1, \dots, X_p)^\top$ be a p -dimensional *nonstochastic* (i.e., fixed) vector.
2. **Data:** Let $\{(\mathbf{x}_i^\top, \mathbf{y}_i^\top)^\top\}_{i=1}^n$ be n i.i.d observations from $(X^\top, Y^\top)^\top$. Define the matrices \mathbf{X} and \mathbf{Y} by

$$\mathbf{X} = [\mathbf{x}_1 \ \cdots \ \mathbf{x}_n] \in \mathbb{R}^{p \times n}, \quad \text{and} \quad \mathbf{Y} = [\mathbf{y}_1 \ \cdots \ \mathbf{y}_n] \in \mathbb{R}^{s \times n},$$

respectively.

Define the following quantities

$$\begin{aligned} \bar{\mathbf{x}} &:= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \\ \bar{\mathbf{y}} &:= \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i, \\ \bar{\mathbf{X}} &:= [\bar{\mathbf{x}} \ \cdots \ \bar{\mathbf{x}}] \in \mathbb{R}^{p \times n}, \\ \bar{\mathbf{Y}} &:= [\bar{\mathbf{y}} \ \cdots \ \bar{\mathbf{y}}] \in \mathbb{R}^{s \times n}. \end{aligned}$$

Defined the centered versions of \mathbf{X} and \mathbf{Y} to be

$$\mathbf{X}_c := \mathbf{X} - \bar{\mathbf{X}}, \quad \text{and} \quad \mathbf{Y}_c := \mathbf{Y} - \bar{\mathbf{Y}}.$$

3. Classical Multivariate Regression Model:

- (a) *Model specification:* Suppose each component of Y depends on the same set of predictors X_1, \dots, X_p in the following way

$$Y_j = \mu_j + \theta_{\theta,1}X_1 + \theta_{\theta,2}X_2 + \cdots + \theta_{\theta,p}X_p + \varepsilon_j, \quad \text{for all } j = 1, \dots, s, \quad (11)$$

where μ_j is the intercept term and ε_j is the error variable with zero mean.

With the data, we can write (11) collectively as

$$\mathbf{Y} = \boldsymbol{\mu} + \boldsymbol{\Theta}\mathbf{X} + \mathbf{E}, \quad (12)$$

where $\boldsymbol{\mu} \in \mathbb{R}^{s \times n}$ is the matrix of the intercept terms, $\boldsymbol{\Theta} \in \mathbb{R}^{s \times p}$ is the matrix of the regression coefficients, and $\mathbf{E} = [\mathbf{E}_1, \mathbf{E}_2, \dots, \mathbf{E}_n]$ is the $s \times n$ error matrix. We assume that each column of \mathbf{E} has mean $\mathbf{0}_s$ and the common unknown nonsingular $s \times s$ covariance matrix $\boldsymbol{\Sigma}_{\mathbf{E}\mathbf{E}}$.

- (b) *Centered version:* We remove $\boldsymbol{\mu}$ from the equation by centering \mathbf{X} and \mathbf{Y} and then estimate $\boldsymbol{\Theta}$ directly. To this end, we let

$$\boldsymbol{\mu} = \bar{\mathbf{Y}} - \boldsymbol{\Theta}\bar{\mathbf{X}}. \quad (13)$$

Then, the model (12) becomes

$$\mathbf{Y}_c = \boldsymbol{\Theta}\mathbf{X}_c + \mathbf{E}. \quad (14)$$

- (c) *Least squares estimation:* Apply the “vec” operation to both sides of (14), we obtain

$$\text{vec}(\mathbf{Y}_c) = (\mathbf{I}_s \otimes \mathbf{X}_c^\top) \text{vec}(\boldsymbol{\Theta}) + \text{vec}(\mathbf{E}), \quad (15)$$

where we have $\text{vec}(\mathbf{Y}_c) \in \mathbb{R}^{sn \times 1}$, $\mathbf{I}_s \otimes \mathbf{X}_c^\top \in \mathbb{R}^{sn \times sp}$, $\text{vec}(\boldsymbol{\Theta}) \in \mathbb{R}^{sp \times 1}$ and $\text{vec}(\mathbf{E}) \in \mathbb{R}^{sn \times 1}$.

Inspecting (15), we see it is a multiple linear regression problem. The error $\text{vec}(\mathbf{E})$ has the mean vector $\mathbf{0}_{sn}$ and $sn \times sn$ block-diagonal covariance matrix

$$\text{Var}[\text{vec}(\mathbf{E})] = \mathbb{E}[\text{vec}(\mathbf{E})\text{vec}(\mathbf{E})^\top] = \boldsymbol{\Sigma}_{\mathbf{E}\mathbf{E}} \otimes \mathbf{I}_n.$$

Assuming $\mathbf{X}_c \mathbf{X}_c^\top$ is nonsingular, we estimate $\text{vec}(\boldsymbol{\Theta})$ using the generalized least-squares method and solve the following minimization problem

$$\underset{\boldsymbol{\Theta}}{\text{minimize}} \left\{ (\text{vec}(\mathbf{Y}_c) - (\mathbf{I}_s \otimes \mathbf{X}_c^\top) \text{vec}(\boldsymbol{\Theta})) [\text{Var}[\text{vec}(\mathbf{E})]]^{-1} (\text{vec}(\mathbf{Y}_c) - (\mathbf{I}_s \otimes \mathbf{X}_c^\top) \text{vec}(\boldsymbol{\Theta})) \right\}.$$

Then,

$$\begin{aligned} \text{vec}(\hat{\boldsymbol{\Theta}}) &= ((\mathbf{I}_s \otimes \mathbf{X}_c)(\boldsymbol{\Sigma}_{\mathbf{E}\mathbf{E}} \otimes \mathbf{I}_n)^{-1}(\mathbf{I}_s \otimes \mathbf{X}_c^\top))^{-1}(\mathbf{I}_s \otimes \mathbf{X}_c)(\boldsymbol{\Sigma}_{\mathbf{E}\mathbf{E}} \otimes \mathbf{I}_n)^{-1} \text{vec}(\mathbf{Y}_c) \\ &= (\mathbf{I}_s \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c) \text{vec}(\mathbf{Y}_c), \end{aligned}$$

using the properties of the Kronecker product.

If we “un-vec” everything, we obtain

$$\hat{\boldsymbol{\Theta}} = \mathbf{Y}_c \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}, \quad (16)$$

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{Y}} - \hat{\boldsymbol{\Theta}} \bar{\mathbf{X}}. \quad (17)$$

- (d) *Minimum-variance linear unbiased estimator of $\text{trace}(\mathbf{A}\boldsymbol{\Theta})$:* Let \mathbf{A} be a fixed matrix. Then, under the conditions above and if $\mathbf{X}_c \mathbf{X}_c^\top$ is nonsingular, then the minimum-variance linear unbiased estimator of $\text{trace}(\mathbf{A}\boldsymbol{\Theta})$ is given by $\text{trace}(\mathbf{A}\hat{\boldsymbol{\Theta}})$.
- (e) *Interpretation of $\hat{\boldsymbol{\Theta}}$:* Suppose we transpose both sides of (14) so that

$$\mathbf{Y}_c^\top = \mathbf{X}_c^\top \boldsymbol{\Theta}^\top + \mathbf{E}^\top.$$

Let $\tilde{\mathbf{y}}_j \in \mathbb{R}^n$ be the j -th column vector of \mathbf{Y}_c^\top , which represents all the n mean-centered observations on the j -th output variable, for $j = 1, \dots, s$. Then, $\tilde{\mathbf{y}}_j \in \mathbb{R}^n$ can be modeled by the multiple regression equation

$$\tilde{\mathbf{y}}_j = \mathbf{X}_c^\top \boldsymbol{\theta}_j + \mathbf{e}_j,$$

where $\boldsymbol{\theta}_j$ is the j -th column of $\boldsymbol{\Theta}^\top$ and \mathbf{e}_j is the j -th column of \mathbf{E}^\top . The ordinary least-square estimator of $\boldsymbol{\theta}_j$ is

$$\hat{\boldsymbol{\theta}}_j = (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c \tilde{\mathbf{y}}_j,$$

which is exactly the j -th row of (16).

Thus, simultaneous (unrestricted) least-squares estimation applied to all the s equations of the multivariate regression model yields the *same* results as does equation-by-equation least-squares. As a result, nothing is gained by estimating the equations jointly, even though the output variables Y may be correlated.

In other words, even though the variables in Y may be correlated, the LS estimator, $\hat{\Theta}$, of Θ does *not* contain any reference to that correlation.

- (f) *Covariance Matrix of $\hat{\Theta}$* : We derive the covariance matrix of $\hat{\Theta}$. By the relationship $\mathbf{Y}_c = \Theta \mathbf{X}_c + \mathbf{E}$, we have

$$\begin{aligned}\hat{\Theta} &= \mathbf{Y}_c \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \\ &= (\Theta \mathbf{X}_c + \mathbf{E}) \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \\ &= \Theta + \mathbf{E} \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}.\end{aligned}$$

It follows that

$$\text{vec}(\hat{\Theta} - \Theta) = \text{vec}(\mathbf{E} \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}) = (\mathbf{I}_s \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c) \text{vec}(\mathbf{E}),$$

and, hence,

$$\begin{aligned}\text{Var}[\text{vec}(\hat{\Theta})] &= \mathbb{E}[\text{vec}(\hat{\Theta} - \Theta) \text{vec}(\hat{\Theta} - \Theta)^\top] \\ &= (\mathbf{I}_s \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c) (\Sigma_{\mathbf{E}\mathbf{E}} \otimes \mathbf{I}_n) (\mathbf{I}_s \otimes \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}) \\ &= \Sigma_{\mathbf{E}\mathbf{E}} \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}.\end{aligned}$$

- (g) *Distribution of $\hat{\Theta}$* : If we assume that the errors in the model (12) are distributed as i.i.d Gaussian random vectors,

$$\mathbf{E}_i \sim \text{Normal}(\mathbf{0}_s, \Sigma_{\mathbf{E}\mathbf{E}}), \quad \text{for all } i = 1, \dots, n,$$

where \mathbf{E}_i denotes the i -th column of \mathbf{E} , then,

$$\text{vec}(\hat{\Theta}) \sim \text{Normal}(\text{vec}(\Theta), \Sigma_{\mathbf{E}\mathbf{E}} \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}).$$

- (h) *Fitted Values*: The $s \times n$ matrix $\hat{\mathbf{Y}}$ of fitted values is given by

$$\hat{\mathbf{Y}} = \hat{\mu} + \hat{\Theta} \mathbf{X} = \bar{\mathbf{Y}} + \hat{\Theta}(\mathbf{X} - \bar{\mathbf{X}}).$$

Also, we have

$$\hat{\mathbf{Y}}_c = \hat{\Theta} \mathbf{X}_c = \mathbf{Y}_c \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c = \mathbf{Y}_c \mathbf{H},$$

where the $n \times n$ matrix $\mathbf{H} = \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c$ is the *hat matrix*.

- (i) *Residual Matrix*: The $s \times n$ residual matrix $\hat{\mathbf{E}}$ is the difference between the observed and fitted values of \mathbf{Y} , i.e.,

$$\hat{\mathbf{E}} := \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y}_c - \hat{\Theta} \mathbf{X}_c = \mathbf{Y}_c (\mathbf{I}_n - \mathbf{H}),$$

In addition, we also have

$$\widehat{\mathbf{E}} = (\boldsymbol{\Theta}\mathbf{X}_c + \mathbf{E}) - (\boldsymbol{\Theta} + \mathbf{E}\mathbf{X}_c^\top(\mathbf{X}_c\mathbf{X}_c^\top)^{-1})\mathbf{X}_c = \mathbf{E}(\mathbf{I}_n - \mathbf{H}).$$

It follows that

$$\begin{aligned}\mathbb{E}[\text{vec}(\widehat{\mathbf{E}})] &= \mathbf{0}_{ns}, \\ \text{Var}[\text{vec}(\widehat{\mathbf{E}})] &= \boldsymbol{\Sigma}_{\mathbf{EE}} \otimes (\mathbf{I}_n - \mathbf{H}).\end{aligned}$$

(j) *Estimation of $\boldsymbol{\Sigma}_{\mathbf{EE}}$* : The $s \times s$ matrix version of the residual sum of squares is

$$\mathbf{S}_{\mathbf{E}} := \widehat{\mathbf{E}}\widehat{\mathbf{E}}^\top = (\mathbf{Y}_c(\mathbf{I}_n - \mathbf{H}))(\mathbf{Y}_c(\mathbf{I}_n - \mathbf{H}))^\top = \mathbf{Y}_c(\mathbf{I}_n - \mathbf{H})\mathbf{Y}_c^\top.$$

Also, it is easy to show

$$\mathbf{S}_{\mathbf{E}} = \mathbf{E}(\mathbf{I}_n - \mathbf{H})\mathbf{E}^\top.$$

Let \mathbf{E}_j be the j -th row of \mathbf{E} . Then, the (j, k) -th element of $\mathbf{S}_{\mathbf{E}}$ can be written as

$$[\mathbf{S}_{\mathbf{E}}]_{(j,k)} = \mathbf{E}_j(\mathbf{I}_n - \mathbf{H})\mathbf{E}_k^\top,$$

whence,

$$\begin{aligned}\mathbb{E}[[\mathbf{S}_{\mathbf{E}}]_{(j,k)}] &= \mathbb{E}[\text{trace}(\mathbf{I}_n - \mathbf{H})\mathbf{E}_k^\top\mathbf{E}_j] \\ &= \text{trace}(\mathbf{I}_n - \mathbf{H})[\boldsymbol{\Sigma}_{\mathbf{EE}}]_{(j,k)} \\ &= (n - p)[\boldsymbol{\Sigma}_{\mathbf{EE}}]_{(j,k)}.\end{aligned}$$

The matrix

$$\widehat{\boldsymbol{\Sigma}}_{\mathbf{EE}} := \frac{1}{n - p}\mathbf{S}_{\mathbf{E}}$$

is called the *residual covariance matrix*.

(k) *Properties of $\widehat{\boldsymbol{\Sigma}}_{\mathbf{EE}}$* :

- i. The matrix $\widehat{\boldsymbol{\Sigma}}_{\mathbf{EE}}$ is statistically independent of $\widehat{\boldsymbol{\Theta}}$;
- ii. The matrix $\widehat{\boldsymbol{\Sigma}}_{\mathbf{EE}}$ has a Wishart distribution with $n - p$ degrees of freedom and expectation $\boldsymbol{\Sigma}_{\mathbf{EE}}$.

(l) *Estimator of $\text{Var}[\widehat{\boldsymbol{\Theta}}]$* : Using the results above, we can estimate $\text{Var}[\widehat{\boldsymbol{\Theta}}] = \boldsymbol{\Sigma}_{\mathbf{EE}} \otimes (\mathbf{X}_c\mathbf{X}_c^\top)^{-1}$ by

$$\widehat{\text{Var}}[\text{vec}(\widehat{\boldsymbol{\Theta}})] = \widehat{\boldsymbol{\Sigma}}_{\mathbf{EE}} \otimes (\mathbf{X}_c\mathbf{X}_c^\top)^{-1}. \quad (18)$$

(m) *Confidence interval*: Let $\boldsymbol{\gamma} \in \mathbb{R}^{sp}$ be an arbitrary vector and consider to construct a confidence interval of $\boldsymbol{\gamma}^\top \text{vec}(\boldsymbol{\Theta})$. Assuming the error vectors are distributed as

$$\mathbf{E}_i \sim \text{Normal}(\mathbf{0}_s, \boldsymbol{\Sigma}_{\mathbf{EE}}), \quad \text{for all } i = 1, \dots, n,$$

we have the quantity

$$t = \frac{\boldsymbol{\gamma}^\top \text{vec}(\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta})}{\sqrt{\boldsymbol{\gamma}^\top (\hat{\boldsymbol{\Sigma}}_{\mathbf{EE}} \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}) \boldsymbol{\gamma}}} \quad (19)$$

has the Student's t -distribution with $n - p$ degrees of freedom. Thus, a $(1 - \alpha) \times 100\%$ confidence interval for $\boldsymbol{\gamma}^\top \text{vec}(\boldsymbol{\Theta})$ can be given by

$$\boldsymbol{\gamma}^\top \text{vec}(\hat{\boldsymbol{\Theta}}) \pm t_{n-p, 1-\frac{\alpha}{2}} \sqrt{\boldsymbol{\gamma}^\top (\hat{\boldsymbol{\Sigma}}_{\mathbf{EE}} \otimes (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}) \boldsymbol{\gamma}},$$

where $t_{n-p, 1-\frac{\alpha}{2}}$ is the $(1 - \alpha/2) \times 100\%$ percentile of the t -distribution with degrees of freedom $n - p$.

- 4. Linear Constrained Estimation:** Consider the following model with centered data matrices \mathbf{X}_c and \mathbf{Y}_c

$$\mathbf{Y}_c = \boldsymbol{\Theta} \mathbf{X}_c + \mathbf{E}.$$

We require $\boldsymbol{\Theta}$ to satisfy a set of known linear constraints of the form

$$\mathbf{K} \boldsymbol{\Theta} \mathbf{L} = \boldsymbol{\Gamma},$$

where the matrix $\mathbf{K} \in \mathbb{R}^{m \times s}$ and the matrix $\mathbf{L} \in \mathbb{R}^{p \times u}$ are full-rank matrices of known constants, and $\boldsymbol{\Gamma} \in \mathbb{R}^{m \times u}$ is a matrix of parameters (known or unknown). We often take $\boldsymbol{\Gamma} = \mathbf{0}_{m \times u}$. We require $m \leq s$ and $u \leq p$.

- (a) *Example — variable selection:* Suppose we wish to study whether a specific subset of the p input variables has little or no effect on the behavior of the output variables. Suppose we arrange the rows of \mathbf{X}_c so that

$$\mathbf{X}_c = \begin{pmatrix} \mathbf{X}_{c,1} \\ \mathbf{X}_{c,2} \end{pmatrix},$$

where $\mathbf{X}_{c,1} \in \mathbb{R}^{p_1 \times n}$ and $\mathbf{X}_{c,2} \in \mathbb{R}^{p_2 \times n}$ with $p_1 + p_2 = p$. Suppose we believe that the variables included in $\mathbf{X}_{c,2}$ do *not* belong in the regression. Corresponding to the partition of \mathbf{X}_c , we set $\boldsymbol{\Theta} = (\boldsymbol{\Theta}_1, \boldsymbol{\Theta}_2)$, so that

$$\mathbf{Y}_c = \boldsymbol{\Theta}_1 \mathbf{X}_{c,1} + \boldsymbol{\Theta}_2 \mathbf{X}_{c,2} + \mathbf{E},$$

where $\boldsymbol{\Theta}_1 \in \mathbb{R}^{s \times p_1}$ and $\boldsymbol{\Theta}_2 \in \mathbb{R}^{s \times p_2}$.

To study whether the input variables included in $\mathbf{X}_{c,2}$ can be eliminated from the model, we set

$$\mathbf{K} = \mathbf{I}_s, \quad \mathbf{L} = \begin{pmatrix} \mathbf{0}_{p_1 \times p_2} \\ \mathbf{I}_{p_2 \times p_2} \end{pmatrix},$$

so that $\mathbf{K} \boldsymbol{\Theta} \mathbf{L} = \boldsymbol{\Theta}_2 = \mathbf{0}_{s \times p_2}$.

- (b) *Constrained least-squares estimation*: To estimate Θ under the linear constraint $\mathbf{K}\Theta\mathbf{L} = \Gamma$, we consider the following optimization problem

$$\begin{aligned} & \underset{\Theta}{\text{minimize}} \left\{ \text{trace} \left((\mathbf{Y}_c - \Theta \mathbf{X}_c)(\mathbf{Y}_c - \Theta \mathbf{X}_c)^\top \right) \right\} \\ & \text{subject to } \mathbf{K}\Theta\mathbf{L} = \Gamma. \end{aligned} \quad (20)$$

Let $\hat{\Theta}^*$ be the minimizer of (20) and Λ be a matrix of Lagrangian coefficients. The normal equations are

$$\begin{aligned} \hat{\Theta}^* \mathbf{X}_c \mathbf{X}_c^\top + \mathbf{K}^\top \Lambda \mathbf{L}^\top &= \mathbf{Y}_c \mathbf{X}_c^\top, \\ \mathbf{K} \hat{\Theta}^* \mathbf{L} &= \Gamma. \end{aligned} \quad (21)$$

It follows that

$$\hat{\Theta}^* = \hat{\Theta} - \mathbf{K}^\top \Lambda \mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1},$$

where $\hat{\Theta} = \mathbf{Y}_c \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}$ as we have shown before. Then,

$$\mathbf{K}(\hat{\Theta} - \mathbf{K}^\top \Lambda \mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1})\mathbf{L} = \Gamma,$$

from which we obtain

$$\Lambda = (\mathbf{K}\mathbf{K}^\top)^{-1}(\mathbf{K}\hat{\Theta}\mathbf{L} - \Gamma)(\mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{L})^{-1},$$

assuming all inverses exist. Finally, we obtain

$$\hat{\Theta}^* = \hat{\Theta} - \mathbf{K}^\top (\mathbf{K}\mathbf{K}^\top)^{-1}(\mathbf{K}\hat{\Theta}\mathbf{L} - \Gamma)(\mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{L})^{-1} \mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1}.$$

- (c) *Multivariate Analysis of Variance (MANOVA)*:

- Residual sum of squares: The residual sum of squares under the constrained model is given by

$$\begin{aligned} \mathbf{S}_E^* &:= (\mathbf{Y}_c - \hat{\Theta}^* \mathbf{X}_c)(\mathbf{Y}_c - \hat{\Theta}^* \mathbf{X}_c)^\top \\ &= (\mathbf{Y}_c - \hat{\Theta} \mathbf{X}_c + (\hat{\Theta} - \hat{\Theta}^*) \mathbf{X}_c)(\mathbf{Y}_c - \hat{\Theta} \mathbf{X}_c + (\hat{\Theta} - \hat{\Theta}^*) \mathbf{X}_c)^\top \\ &= (\mathbf{Y}_c - \hat{\Theta} \mathbf{X}_c)(\mathbf{Y}_c - \hat{\Theta} \mathbf{X}_c)^\top + (\hat{\Theta} - \hat{\Theta}^*) \mathbf{X}_c \mathbf{X}_c^\top (\hat{\Theta} - \hat{\Theta}^*)^\top, \end{aligned} \quad (22)$$

where

- the first term on the RHS of (22) is the matrix version of the residual sum of squares, \mathbf{S}_E , for the *unconstrained* model, and
- the second term is the additional source of variation, $\mathbf{S}_h := \mathbf{S}_E - \mathbf{S}_E^*$, due to dropping the constraints.

- Regression sum of squares: The *regression sum of squares*, \mathbf{S}_{reg} , for the unconstrained model is given by

$$\begin{aligned}\mathbf{S}_{\text{reg}} &:= \widehat{\boldsymbol{\Theta}} \mathbf{X}_c \mathbf{X}_c^\top \widehat{\boldsymbol{\Theta}}^\top \\ &= (\widehat{\boldsymbol{\Theta}}^* + (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*)) \mathbf{X}_c \mathbf{X}_c^\top (\widehat{\boldsymbol{\Theta}}^* + (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*))^\top \\ &= \widehat{\boldsymbol{\Theta}}^* \mathbf{X}_c \mathbf{X}_c^\top \widehat{\boldsymbol{\Theta}}^* + (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*) \mathbf{X}_c \mathbf{X}_c^\top (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*)^\top,\end{aligned}\quad (23)$$

where

- the first term on the RHS of (23) is $\mathbf{S}_{\text{reg}}^*$, the matrix version of the regression sum of squares for the constrained model, and
- the second term is, again, \mathbf{S}_h .
- MANOVA Table: Let $k = \text{rank}(\mathbf{K})$. We have the following MANOVA table.

| Source of Variation | df | Sum of Squares |
|-----------------------------|-------------|--|
| Constrained Model | $p - k$ | $\mathbf{S}_{\text{reg}}^* = \widehat{\boldsymbol{\Theta}}^* \mathbf{X}_c \mathbf{X}_c^\top \widehat{\boldsymbol{\Theta}}^*$ |
| Due to dropping constraints | k | $\mathbf{S}_h = (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*) \mathbf{X}_c \mathbf{X}_c^\top (\widehat{\boldsymbol{\Theta}} - \widehat{\boldsymbol{\Theta}}^*)^\top$ |
| Unconstrained model | p | $\mathbf{S}_{\text{reg}} = \widehat{\boldsymbol{\Theta}} \mathbf{X}_c \mathbf{X}_c^\top \widehat{\boldsymbol{\Theta}}$ |
| Residual | $n - p - 1$ | $\mathbf{S}_E = (\mathbf{Y}_c - \widehat{\boldsymbol{\Theta}} \mathbf{X}_c)(\mathbf{Y}_c - \widehat{\boldsymbol{\Theta}} \mathbf{X}_c)^\top$ |
| Total | $n - 1$ | $\mathbf{Y}_c \mathbf{Y}_c^\top$ |

(d) *Expectation of \mathbf{S}_{reg}* : With the previous results, we have

$$\mathbb{E}[\mathbf{S}_h] = \mathbf{D}(\mathbf{K}\widehat{\boldsymbol{\Theta}}\mathbf{L} - \boldsymbol{\Gamma})(\mathbf{L}^\top(\mathbf{X}_c \mathbf{X}_c^\top)^{-1}\mathbf{L})^{-1}(\mathbf{K}\widehat{\boldsymbol{\Theta}}\mathbf{L} - \boldsymbol{\Gamma})^\top \mathbf{D}^\top + \mathbf{F} \mathbb{E}[\mathbf{E}\mathbf{G}\mathbf{E}^\top] \mathbf{F}^\top,$$

where

$$\begin{aligned}\mathbf{D} &= \mathbf{K}^\top (\mathbf{K}\mathbf{K}^\top)^{-1}, \\ \mathbf{F} &= \mathbf{D}\mathbf{K}, \\ \mathbf{G} &= \mathbf{X}_c^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{L} (\mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{L})^{-1} \mathbf{L}^\top (\mathbf{X}_c \mathbf{X}_c^\top)^{-1} \mathbf{X}_c^\top.\end{aligned}$$

It is easy to show

$$\mathbf{F}^2 = \mathbf{F} = \mathbf{F}^\top, \quad \text{and} \quad \mathbf{G}^2 = \mathbf{G} = \mathbf{G}^\top,$$

and

$$\mathbb{E}[\mathbf{E}\mathbf{G}\mathbf{E}^\top] = u \boldsymbol{\Sigma}_{\mathbf{EE}},$$

where $u = \text{trace}(\mathbf{G}) = \text{trace}(\mathbf{I}_u)$ is the rank of \mathbf{L} .

(e) *Hypothesis testing*: We test

$$H_0 : \mathbf{K}\boldsymbol{\Theta}\mathbf{L} = \boldsymbol{\Gamma} \quad \text{vs.} \quad H_1 : \mathbf{K}\boldsymbol{\Theta}\mathbf{L} \neq \boldsymbol{\Gamma}.$$

Under H_0 ,

$$\mathbb{E} \left[\frac{1}{u} \mathbf{S}_h \right] = \mathbf{F} \boldsymbol{\Sigma}_{\mathbf{EE}} \mathbf{F}^\top,$$

$$\mathbb{E} \left[\frac{1}{n - p - 1} \mathbf{S}_{\mathbf{E}} \right] = \boldsymbol{\Sigma}_{\mathbf{EE}}.$$

A formal significance test of H_0 vs. H_1 can be realized through a function (e.g., determinant, trace, or largest eigenvalue) of the quantity $\mathbf{F} \mathbf{S}_h \mathbf{F}^{-1} (\mathbf{F} \mathbf{S}_{\mathbf{E}} \mathbf{F}^\top)^{-1}$. Examples include

- Hotelling-Lawley trace statistic: $\text{trace}(\mathbf{S}_h \mathbf{S}_{\mathbf{E}}^{-1})$,
- Roy's largest root: $\lambda_{\max}(\mathbf{S}_h \mathbf{S}_{\mathbf{E}}^{-1})$, and
- Wilks's lambda (likelihood ratio criterion): $|\mathbf{S}_{\mathbf{E}}|/|\mathbf{S}_h + \mathbf{S}_{\mathbf{E}}|$.

Under appropriate distribution assumptions, we reject H_0 if Hotelling-Lawley's trace statistic and Roy's largest root are small, or if Wilk's lambda is large.

References

Izenman, Alan J (Mar. 2009). *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*. en. Springer Science & Business Media.