

Chapter 11

Multivariate Regression

11.1 Introduction

Multivariate regression is a system of regression equations. Multivariate regression is used as reduced form models for instrumental variable estimation (Chapter 12), vector autoregressions (Chapter 15), demand systems (demand for multiple goods), and other contexts.

Multivariate regression is also called by the name **systems of regression equations**. Closely related is the method of **Seemingly Unrelated Regressions** (SUR) introduced in Section 11.7.

Most of the tools of single equation regression generalize to multivariate regression. A major difference is a new set of notation to handle matrix estimators.

11.2 Regression Systems

A univariate linear regression equation equals $Y = X'\beta + e$ where Y is scalar and X is a vector. **Multivariate regression** is a system of m linear regressions, and equals

$$Y_j = X_j'\beta_j + e_j \quad (11.1)$$

for $j = 1, \dots, m$. Here we use the subscript j to denote the j^{th} dependent variable, not the i^{th} individual. As an example, Y_j could be expenditures by a household on good category j (e.g., food, housing, transportation, clothing, recreation). The regressor vectors X_j are $k_j \times 1$ and e_j is an error. The coefficient vectors β_j are $k_j \times 1$. The total number of coefficients are $\bar{k} = \sum_{j=1}^m k_j$. The regressors can be common across j or can vary across j . In the household expenditure example the regressors X_j are typically common across j , and include variables such as household income, number and ages of family members, and demographic characteristics. The regression system specializes to univariate regression when $m = 1$.

Define the $m \times 1$ error vector $e = (e_1, \dots, e_m)'$ and its $m \times m$ covariance matrix $\Sigma = \mathbb{E}[ee']$. The diagonal elements are the variances of the errors e_j and the off-diagonals are the covariances across variables.

We can group the m equations (11.1) into a single equation as follows. Let $Y = (Y_1, \dots, Y_m)'$ be the $m \times 1$ vector of dependent variables. Define the $m \times \bar{k}$ matrix of regressors

$$\bar{X} = \begin{pmatrix} X_1' & 0 & \cdots & 0 \\ \vdots & X_2' & & \vdots \\ 0 & 0 & \cdots & X_m' \end{pmatrix}$$

and the $\bar{k} \times 1$ stacked coefficient vector

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}.$$

The m regression equations can be jointly written as

$$Y = \bar{X}\beta + e. \quad (11.2)$$

This is a system of m equations.

For n observations the joint system can be written in matrix notation by stacking. Define

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}, \quad \bar{\mathbf{X}} = \begin{pmatrix} \bar{X}_1 \\ \vdots \\ \bar{X}_n \end{pmatrix}$$

which are $mn \times 1$, $mn \times 1$, and $mn \times \bar{k}$, respectively. The system can be written as $\mathbf{Y} = \bar{\mathbf{X}}\beta + \mathbf{e}$.

In many applications the regressor vectors X_j are common across the variables j , so $X_j = X$ and $k_j = k$. By this we mean that the same variables enter each equation with no exclusion restrictions. Several important simplifications occur in this context. One is that we can write (11.2) using the notation

$$Y = \mathbf{B}'X + e \quad (11.3)$$

where $\mathbf{B} = (\beta_1, \beta_2, \dots, \beta_m)$ is $k \times m$. Another is that we can write the joint system of observations in the $n \times m$ matrix notation $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$ where

$$\mathbf{Y} = \begin{pmatrix} Y'_1 \\ \vdots \\ Y'_n \end{pmatrix}, \quad \mathbf{E} = \begin{pmatrix} e'_1 \\ \vdots \\ e'_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} X'_1 \\ \vdots \\ X'_n \end{pmatrix}.$$

Another convenient implication of common regressors is that we have the simplification

$$\bar{\mathbf{X}} = \begin{pmatrix} X' & 0 & \cdots & 0 \\ 0 & X' & & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & X' \end{pmatrix} = \mathbf{I}_m \otimes X'$$

where \otimes is the Kronecker product (see Appendix A.21).

11.3 Least Squares Estimator

The equations (11.1) can be estimated by least squares. This takes the form

$$\hat{\beta}_j = \left(\sum_{i=1}^n X_{ji} X'_{ji} \right)^{-1} \left(\sum_{i=1}^n X_{ji} Y_{ji} \right).$$

An estimator of β is the stacked vector

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_m \end{pmatrix}.$$

We can alternatively write this estimator using the systems notation

$$\hat{\beta} = (\overline{\mathbf{X}}' \overline{\mathbf{X}})^{-1} (\overline{\mathbf{X}}' \mathbf{Y}) = \left(\sum_{i=1}^n \overline{X}_i' \overline{X}_i \right)^{-1} \left(\sum_{i=1}^n \overline{X}_i' Y_i \right). \quad (11.4)$$

To see this, observe that

$$\begin{aligned} \overline{\mathbf{X}}' \overline{\mathbf{X}} &= \begin{pmatrix} \overline{X}_1' & \cdots & \overline{X}_n' \end{pmatrix} \begin{pmatrix} \overline{X}_1 \\ \vdots \\ \overline{X}_n \end{pmatrix} \\ &= \sum_{i=1}^n \overline{X}_i' \overline{X}_i \\ &= \sum_{i=1}^n \begin{pmatrix} X_{1i} & 0 & \cdots & 0 \\ \vdots & X_{2i} & & \vdots \\ 0 & 0 & \cdots & X_{mi} \end{pmatrix} \begin{pmatrix} X_{1i}' & 0 & \cdots & 0 \\ \vdots & X_{2i}' & & \vdots \\ 0 & 0 & \cdots & X_{mi}' \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^n X_{1i} X_{1i}' & 0 & \cdots & 0 \\ \vdots & \sum_{i=1}^n X_{2i} X_{2i}' & & \vdots \\ 0 & 0 & \cdots & \sum_{i=1}^n X_{mi} X_{mi}' \end{pmatrix} \end{aligned}$$

and

$$\begin{aligned} \overline{\mathbf{X}}' \mathbf{Y} &= \begin{pmatrix} \overline{X}_1' & \cdots & \overline{X}_n' \end{pmatrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \\ &= \sum_{i=1}^n \overline{X}_i' Y_i \\ &= \sum_{i=1}^n \begin{pmatrix} X_{1i} & 0 & \cdots & 0 \\ \vdots & X_{2i} & & \vdots \\ 0 & 0 & \cdots & X_{mi} \end{pmatrix} \begin{pmatrix} Y_{1i} \\ \vdots \\ Y_{mi} \end{pmatrix} \\ &= \begin{pmatrix} \sum_{i=1}^n X_{1i} Y_{1i} \\ \vdots \\ \sum_{i=1}^n X_{mi} Y_{mi} \end{pmatrix}. \end{aligned}$$

Hence

$$\begin{aligned} (\overline{\mathbf{X}}' \overline{\mathbf{X}})^{-1} (\overline{\mathbf{X}}' \mathbf{Y}) &= \left(\sum_{i=1}^n \overline{X}_i \overline{X}_i' \right)^{-1} \left(\sum_{i=1}^n \overline{X}_i Y_i \right) \\ &= \begin{pmatrix} (\sum_{i=1}^n X_{1i} X_{1i}')^{-1} (\sum_{i=1}^n X_{1i} Y_{1i}) \\ \vdots \\ (\sum_{i=1}^n X_{mi} X_{mi}')^{-1} (\sum_{i=1}^n X_{mi} Y_{mi}) \end{pmatrix} \\ &= \hat{\beta} \end{aligned}$$

as claimed.

The $m \times 1$ residual vector for the i^{th} observation is $\hat{e}_i = Y_i - \bar{X}_i' \hat{\beta}$. The least squares estimator of the $m \times m$ error covariance matrix is

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \hat{e}_i \hat{e}_i'. \quad (11.5)$$

In the case of common regressors, the least squares coefficients can be written as

$$\hat{\beta}_j = \left(\sum_{i=1}^n X_i X_i' \right)^{-1} \left(\sum_{i=1}^n X_i Y_{ji} \right)$$

and

$$\hat{B} = (\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_m) = (X'X)^{-1} (X'Y). \quad (11.6)$$

In Stata, multivariate regression can be implemented using the `mvreg` command.

11.4 Expectation and Variance of Systems Least Squares

We can calculate the finite-sample expectation and variance of $\hat{\beta}$ under the conditional expectation assumption

$$\mathbb{E}[e | X] = 0 \quad (11.7)$$

where X is the union of the regressors X_j . Equation (11.7) is equivalent to $\mathbb{E}[Y_j | X] = X_j' \beta_j$, which means that the regression model is correctly specified.

We can center the estimator as

$$\hat{\beta} - \beta = (\bar{X}'\bar{X})^{-1} (\bar{X}'e) = \left(\sum_{i=1}^n \bar{X}_i' \bar{X}_i \right)^{-1} \left(\sum_{i=1}^n \bar{X}_i' e_i \right).$$

Taking conditional expectations we find $\mathbb{E}[\hat{\beta} | X] = \beta$. Consequently, systems least squares is unbiased under correct specification.

To compute the variance of the estimator, define the conditional covariance matrix of the errors of the i^{th} observation $\mathbb{E}[e_i e_i' | X_i] = \Sigma_i$ which in general is a function of X_i . If the observations are mutually independent then

$$\mathbb{E}[ee' | X] = \mathbb{E} \left[\begin{pmatrix} e_1 e_1' & e_1 e_2' & \cdots & e_1 e_n' \\ \vdots & \ddots & & \vdots \\ e_n e_1' & e_n e_2' & \cdots & e_n e_n' \end{pmatrix} \middle| X \right] = \begin{pmatrix} \Sigma_1 & 0 & \cdots & 0 \\ \vdots & \ddots & & \vdots \\ 0 & 0 & \cdots & \Sigma_n \end{pmatrix}.$$

Also, by independence across observations,

$$\text{var} \left[\sum_{i=1}^n \bar{X}_i' e_i \middle| X \right] = \sum_{i=1}^n \text{var} [\bar{X}_i' e_i | X_i] = \sum_{i=1}^n \bar{X}_i' \Sigma_i \bar{X}_i.$$

It follows that

$$\text{var}[\hat{\beta} | X] = (\bar{X}'\bar{X})^{-1} \left(\sum_{i=1}^n \bar{X}_i' \Sigma_i \bar{X}_i \right) (\bar{X}'\bar{X})^{-1}.$$

When the regressors are common so that $\bar{X}_i = I_m \otimes X_i'$ then the covariance matrix can be written as

$$\text{var}[\hat{\beta} | X] = (I_m \otimes (X'X)^{-1}) \left(\sum_{i=1}^n (\Sigma_i \otimes X_i X_i') \right) (I_m \otimes (X'X)^{-1}).$$

If the errors are conditionally homoskedastic

$$\mathbb{E}[ee' | X] = \Sigma \quad (11.8)$$

then the covariance matrix simplifies to

$$\text{var}[\hat{\beta} | X] = (\bar{X}'\bar{X})^{-1} \left(\sum_{i=1}^n \bar{X}_i' \Sigma \bar{X}_i \right) (\bar{X}'\bar{X})^{-1}.$$

If both simplifications (common regressors and conditional homoskedasticity) hold then we have the considerable simplification

$$\text{var}[\hat{\beta} | X] = \Sigma \otimes (X'X)^{-1}.$$

11.5 Asymptotic Distribution

For an asymptotic distribution it is sufficient to consider the equation-by-equation projection model in which case

$$\mathbb{E}[X_j e_j] = 0. \quad (11.9)$$

First, consider consistency. Since $\hat{\beta}_j$ are the standard least squares estimators, they are consistent for the projection coefficients β_j .

Second, consider the asymptotic distribution. Our single equation theory implies that the $\hat{\beta}_j$ are asymptotically normal. But this theory does not provide a joint distribution of the $\hat{\beta}_j$ across j , which we now derive. Since the vector

$$\bar{X}_i' e_i = \begin{pmatrix} X_{1i} e_{1i} \\ \vdots \\ X_{mi} e_{mi} \end{pmatrix}$$

is i.i.d. across i and mean zero under (11.9), the central limit theorem implies

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \bar{X}_i' e_i \xrightarrow{d} N(0, \Omega)$$

where

$$\Omega = \mathbb{E}[\bar{X}_i' e_i e_i' \bar{X}_i] = \mathbb{E}[\bar{X}_i' \Sigma_i \bar{X}_i].$$

The matrix Ω is the covariance matrix of the variables $X_{ji} e_{ji}$ across equations. Under conditional homoskedasticity (11.8) the matrix Ω simplifies to

$$\Omega = \mathbb{E}[\bar{X}_i' \Sigma \bar{X}_i] \quad (11.10)$$

(see Exercise 11.1). When the regressors are common it simplifies to

$$\Omega = \mathbb{E}[ee' \otimes XX'] \quad (11.11)$$

(see Exercise 11.2). Under both conditions (homoskedasticity and common regressors) it simplifies to

$$\Omega = \Sigma \otimes \mathbb{E}[XX'] \quad (11.12)$$

(see Exercise 11.3).

Applied to the centered and normalized estimator we obtain the asymptotic distribution.

Theorem 11.1 Under Assumption 7.2, $\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, V_\beta)$ where $V_\beta = Q^{-1}\Omega Q^{-1}$ and

$$Q = E[\bar{X}'\bar{X}] = \begin{pmatrix} E[X_1 X_1'] & 0 & \cdots & 0 \\ \vdots & \ddots & & \vdots \\ 0 & 0 & \cdots & E[X_m X_m'] \end{pmatrix}.$$

For a proof, see Exercise 11.4.

When the regressors are common the matrix Q simplifies as

$$Q = I_m \otimes E[XX'] \quad (11.13)$$

(See Exercise 11.5).

If both the regressors are common and the errors are conditionally homoskedastic (11.8) then we have the simplification

$$V_\beta = \Sigma \otimes (E[XX'])^{-1} \quad (11.14)$$

(see Exercise 11.6).

Sometimes we are interested in parameters $\theta = r(\beta_1, \dots, \beta_m) = r(\beta)$ which are functions of the coefficients from multiple equations. In this case the least squares estimator of θ is $\hat{\theta} = r(\hat{\beta})$. The asymptotic distribution of $\hat{\theta}$ can be obtained from Theorem 11.1 by the delta method.

Theorem 11.2 Under Assumptions 7.2 and 7.3, $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, V_\theta)$ where $V_\theta = R' V_\beta R$ and $R = \frac{\partial}{\partial \beta} r(\beta)'$.

For a proof, see Exercise 11.7.

Theorem 11.2 is an example where multivariate regression is fundamentally distinct from univariate regression. Only by treating least squares as a joint estimator can we obtain a distributional theory for a function of multiple equations. We can thereby construct standard errors, confidence intervals, and hypothesis tests.

11.6 Covariance Matrix Estimation

From the finite sample and asymptotic theory we can construct appropriate estimators for the variance of $\hat{\beta}$. In the general case we have

$$\hat{V}_{\hat{\beta}} = (\bar{X}'\bar{X})^{-1} \left(\sum_{i=1}^n \bar{X}_i' \hat{e}_i \hat{e}_i' \bar{X}_i \right) (\bar{X}'\bar{X})^{-1}.$$

Under conditional homoskedasticity (11.8) an appropriate estimator is

$$\hat{V}_{\hat{\beta}}^0 = (\bar{X}'\bar{X})^{-1} \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma} \bar{X}_i \right) (\bar{X}'\bar{X})^{-1}.$$

When the regressors are common then these estimators equal

$$\hat{\mathbf{V}}_{\hat{\beta}} = \left(\mathbf{I}_m \otimes (\mathbf{X}'\mathbf{X})^{-1} \right) \left(\sum_{i=1}^n (\hat{e}_i \hat{e}_i' \otimes X_i X_i') \right) \left(\mathbf{I}_m \otimes (\mathbf{X}'\mathbf{X})^{-1} \right)$$

and $\hat{\mathbf{V}}_{\hat{\beta}}^0 = \hat{\Sigma} \otimes (\mathbf{X}'\mathbf{X})^{-1}$, respectively.

Covariance matrix estimators for $\hat{\theta}$ are found as

$$\begin{aligned} \hat{\mathbf{V}}_{\hat{\theta}} &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}} \hat{\mathbf{R}} \\ \hat{\mathbf{V}}_{\hat{\theta}}^0 &= \hat{\mathbf{R}}' \hat{\mathbf{V}}_{\hat{\beta}}^0 \hat{\mathbf{R}} \\ \hat{\mathbf{R}} &= \frac{\partial}{\partial \beta} r(\beta)'. \end{aligned}$$

Theorem 11.3 Under Assumption 7.2, $n\hat{\mathbf{V}}_{\hat{\beta}} \xrightarrow{p} \mathbf{V}_{\beta}$ and $n\hat{\mathbf{V}}_{\hat{\beta}}^0 \xrightarrow{p} \mathbf{V}_{\beta}^0$.

For a proof, see Exercise 11.8.

11.7 Seemingly Unrelated Regression

Consider the systems regression model under the conditional expectation and homoskedasticity assumptions

$$\begin{aligned} Y &= \bar{X}\beta + e \\ \mathbb{E}[e | X] &= 0 \\ \mathbb{E}[ee' | X] &= \Sigma. \end{aligned} \tag{11.15}$$

Since the errors are correlated across equations we consider estimation by Generalized Least Squares (GLS). To derive the estimator, premultiply (11.15) by $\Sigma^{-1/2}$ so that the transformed error vector is i.i.d. with covariance matrix \mathbf{I}_m . Then apply least squares and rearrange to find

$$\hat{\beta}_{\text{gls}} = \left(\sum_{i=1}^n \bar{X}_i' \Sigma^{-1} \bar{X}_i \right)^{-1} \left(\sum_{i=1}^n \bar{X}_i' \Sigma^{-1} Y_i \right). \tag{11.16}$$

(see Exercise 11.9). Another approach is to take the vector representation

$$\mathbf{Y} = \bar{\mathbf{X}}\beta + \mathbf{e}$$

and calculate that the equation error \mathbf{e} has variance $\mathbb{E}[\mathbf{e}\mathbf{e}'] = \mathbf{I}_n \otimes \Sigma$. Premultiply the equation by $\mathbf{I}_n \otimes \Sigma^{-1/2}$ so that the transformed error has covariance matrix \mathbf{I}_{nm} and then apply least squares to find

$$\hat{\beta}_{\text{gls}} = \left(\bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \bar{\mathbf{X}} \right)^{-1} \left(\bar{\mathbf{X}}' (\mathbf{I}_n \otimes \Sigma^{-1}) \mathbf{Y} \right) \tag{11.17}$$

(see Exercise 11.10).

Expressions (11.16) and (11.17) are algebraically equivalent. To see the equivalence, observe that

$$\begin{aligned}\bar{\mathbf{X}}'(\mathbf{I}_n \otimes \Sigma^{-1})\bar{\mathbf{X}} &= \begin{pmatrix} \bar{\mathbf{X}}_1' & \cdots & \bar{\mathbf{X}}_n' \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & 0 & \cdots & 0 \\ \vdots & \Sigma^{-1} & & \vdots \\ 0 & 0 & \cdots & \Sigma^{-1} \end{pmatrix} \begin{pmatrix} \bar{\mathbf{X}}_1 \\ \vdots \\ \bar{\mathbf{X}}_n \end{pmatrix} \\ &= \sum_{i=1}^n \bar{\mathbf{X}}_i' \Sigma^{-1} \bar{\mathbf{X}}_i\end{aligned}$$

and

$$\begin{aligned}\bar{\mathbf{X}}'(\mathbf{I}_n \otimes \Sigma^{-1})\mathbf{Y} &= \begin{pmatrix} \bar{\mathbf{X}}_1' & \cdots & \bar{\mathbf{X}}_n' \end{pmatrix} \begin{pmatrix} \Sigma^{-1} & 0 & \cdots & 0 \\ \vdots & \Sigma^{-1} & & \vdots \\ 0 & 0 & \cdots & \Sigma^{-1} \end{pmatrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \\ &= \sum_{i=1}^n \bar{\mathbf{X}}_i' \Sigma^{-1} Y_i.\end{aligned}$$

Since Σ is unknown it must be replaced by an estimator. Using $\hat{\Sigma}$ from (11.5) we obtain a feasible GLS estimator.

$$\begin{aligned}\hat{\beta}_{\text{sur}} &= \left(\sum_{i=1}^n \bar{\mathbf{X}}_i' \hat{\Sigma}^{-1} \bar{\mathbf{X}}_i \right)^{-1} \left(\sum_{i=1}^n \bar{\mathbf{X}}_i' \hat{\Sigma}^{-1} Y_i \right) \\ &= \left(\bar{\mathbf{X}}'(\mathbf{I}_n \otimes \hat{\Sigma}^{-1})\bar{\mathbf{X}} \right)^{-1} \left(\bar{\mathbf{X}}'(\mathbf{I}_n \otimes \hat{\Sigma}^{-1})\mathbf{Y} \right).\end{aligned}\tag{11.18}$$

This is the **Seemingly Unrelated Regression (SUR)** estimator as introduced by Zellner (1962).

The estimator $\hat{\Sigma}$ can be updated by calculating the SUR residuals $\hat{e}_i = Y_i - \bar{\mathbf{X}}_i' \hat{\beta}_{\text{sur}}$ and the covariance matrix estimator $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \hat{e}_i \hat{e}_i'$. Substituted into (11.18) we obtain an iterated SUR estimator. This can be iterated until convergence.

Under conditional homoskedasticity (11.8) we can derive its asymptotic distribution.

Theorem 11.4 Under Assumption 7.2 and (11.8)

$$\sqrt{n}(\hat{\beta}_{\text{sur}} - \beta) \xrightarrow{d} N(0, \mathbf{V}_{\beta}^*)$$

where $\mathbf{V}_{\beta}^* = \left(\mathbb{E} \left[\bar{\mathbf{X}}' \Sigma^{-1} \bar{\mathbf{X}} \right] \right)^{-1}$.

For a proof, see Exercise 11.11.

Under these assumptions, SUR is more efficient than least squares.

Theorem 11.5 Under Assumption 7.2 and (11.8)

$$\mathbf{V}_{\beta}^* = \left(\mathbb{E} \left[\bar{\mathbf{X}}' \Sigma^{-1} \bar{\mathbf{X}} \right] \right)^{-1} \leq \left(\mathbb{E} \left[\bar{\mathbf{X}}' \bar{\mathbf{X}} \right] \right)^{-1} \mathbb{E} \left[\bar{\mathbf{X}}' \Sigma \bar{\mathbf{X}} \right] \left(\mathbb{E} \left[\bar{\mathbf{X}}' \bar{\mathbf{X}} \right] \right)^{-1} = \mathbf{V}_{\beta}$$

and thus $\hat{\beta}_{\text{sur}}$ is asymptotically more efficient than $\hat{\beta}_{\text{ols}}$.

For a proof, see Exercise 11.12.

An appropriate estimator of the variance of $\hat{\beta}_{\text{sur}}$ is

$$\hat{V}_{\hat{\beta}} = \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma}^{-1} \bar{X}_i \right)^{-1}.$$

Theorem 11.6 Under Assumption 7.2 and (11.8) $n\hat{V}_{\hat{\beta}} \xrightarrow{p} V_{\beta}$.

For a proof, see Exercise 11.13.

In Stata, the seemingly unrelated regressions estimator is implemented using the `sureg` command.

Arnold Zellner

Arnold Zellner (1927-2010) of the United States was a founding father of the econometrics field. He was a pioneer in Bayesian econometrics. One of his core contributions was the method of Seemingly Unrelated Regressions.

11.8 Equivalence of SUR and Least Squares

When the regressors are common across equations $X_j = X$ it turns out that the SUR estimator simplifies to least squares.

To see this, recall that when regressors are common this implies that $\bar{X} = \mathbf{I}_m \otimes X'$. Then

$$\begin{aligned} \bar{X}_i' \hat{\Sigma}^{-1} &= (\mathbf{I}_m \otimes X_i) \hat{\Sigma}^{-1} \\ &= \hat{\Sigma}^{-1} \otimes X_i \\ &= (\hat{\Sigma}^{-1} \otimes \mathbf{I}_k) (\mathbf{I}_m \otimes X_i) \\ &= (\hat{\Sigma}^{-1} \otimes \mathbf{I}_k) \bar{X}_i'. \end{aligned}$$

Thus

$$\begin{aligned} \hat{\beta}_{\text{sur}} &= \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma}^{-1} \bar{X}_i \right)^{-1} \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma}^{-1} Y_i \right) \\ &= \left((\hat{\Sigma}^{-1} \otimes \mathbf{I}_k) \sum_{i=1}^n \bar{X}_i' \bar{X}_i \right)^{-1} \left((\hat{\Sigma}^{-1} \otimes \mathbf{I}_k) \sum_{i=1}^n \bar{X}_i' Y_i \right) \\ &= \left(\sum_{i=1}^n \bar{X}_i' \bar{X}_i \right)^{-1} \left(\sum_{i=1}^n \bar{X}_i' Y_i \right) = \hat{\beta}_{\text{ols}}. \end{aligned}$$

A model where regressors are not common across equations is nested within a model with the union of all regressors included in all equations. Thus the model with regressors common across equations

is a fully unrestricted model, and a model where the regressors differ across equations is a restricted model. Thus the above result shows that the SUR estimator reduces to least squares in the absence of restrictions, but SUR can differ from least squares otherwise.

Another context where SUR=OLS is when the variance matrix is diagonal, $\Sigma = \text{diag}\{\sigma_1^2, \dots, \sigma_m^2\}$. In this case $\Sigma^{-1/2}\bar{X}_i = \bar{X}_i \text{diag}\{\mathbf{I}_{k_1}\sigma_1^{-1/2}, \dots, \mathbf{I}_{k_m}\sigma_m^{-1/2}\}$ from which you can calculate that $\hat{\beta}_{\text{sur}} = \hat{\beta}_{\text{ols}}$. The intuition is that there is no difference in systems estimation when the equations are uncorrelated, which occurs when Σ is diagonal.

11.9 Maximum Likelihood Estimator

Take the linear model under the assumption that the error is independent of the regressors and multivariate normally distributed. Thus $Y = \bar{X}\beta + e$ with $e \sim N(0, \Sigma)$. In this case we can consider the maximum likelihood estimator (MLE) of the coefficients.

It is convenient to reparameterize the covariance matrix in terms of its inverse $\mathbf{S} = \Sigma^{-1}$. With this reparameterization the conditional density of Y given $X = x$ equals

$$f(y|x) = \frac{\det(\mathbf{S})^{1/2}}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}(y-x\beta)' \mathbf{S}(y-x\beta)\right).$$

The log-likelihood function for the sample is

$$\ell_n(\beta, \mathbf{S}) = -\frac{nm}{2} \log(2\pi) + \frac{n}{2} \log(\det(\mathbf{S})) - \frac{1}{2} \sum_{i=1}^n (Y_i - \bar{X}_i\beta)' \mathbf{S} (Y_i - \bar{X}_i\beta).$$

The maximum likelihood estimator $(\hat{\beta}_{\text{mle}}, \hat{\mathbf{S}}_{\text{mle}})$ maximizes the log-likelihood function. The first order conditions are

$$0 = \frac{\partial}{\partial \beta} \ell_n(\beta, \mathbf{S}) \Big|_{\beta=\hat{\beta}, \mathbf{S}=\hat{\mathbf{S}}} = \sum_{i=1}^n \bar{X}_i \hat{\mathbf{S}} (Y_i - \bar{X}_i \hat{\beta})$$

and

$$0 = \frac{\partial}{\partial \mathbf{S}} \ell_n(\beta, \mathbf{S}) \Big|_{\beta=\hat{\beta}, \mathbf{S}=\hat{\mathbf{S}}} = \frac{n}{2} \hat{\mathbf{S}}^{-1} - \frac{1}{2} \text{tr} \left(\sum_{i=1}^n (Y_i - \bar{X}_i \hat{\beta})(Y_i - \bar{X}_i \hat{\beta})' \right).$$

The second equation uses the matrix results $\frac{\partial}{\partial \mathbf{S}} \log(\det(\mathbf{S})) = \mathbf{S}^{-1}$ and $\frac{\partial}{\partial \mathbf{B}} \text{tr}(\mathbf{AB}) = \mathbf{A}'$ from Appendix A.20.

Solving and making the substitution $\hat{\Sigma} = \hat{\mathbf{S}}^{-1}$ we obtain

$$\hat{\beta}_{\text{mle}} = \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma}^{-1} \bar{X}_i \right)^{-1} \left(\sum_{i=1}^n \bar{X}_i' \hat{\Sigma}^{-1} Y_i \right)$$

$$\hat{\Sigma}_{\text{mle}} = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{X}_i \hat{\beta})(Y_i - \bar{X}_i \hat{\beta})'.$$

Notice that each equation refers to the other. Hence these are not closed-form expressions but can be solved via iteration. The solution is identical to the iterated SUR estimator. Thus the iterated SUR estimator is identical to MLE under normality.

Recall that the SUR estimator simplifies to OLS when the regressors are common across equations. The same occurs for the MLE. Thus when $\bar{X}_i = \mathbf{I}_m \otimes X_i'$ we find that $\hat{\beta}_{\text{mle}} = \hat{\beta}_{\text{ols}}$ and $\hat{\Sigma}_{\text{mle}} = \hat{\Sigma}_{\text{ols}}$.

11.10 Restricted Estimation

In many multivariate regression applications it is desired to impose restrictions on the coefficients. In particular, cross-equation restrictions (for example, imposing Slutsky symmetry on a demand system) can be quite important and can only be imposed by a multivariate estimation method. Estimation subject to restrictions can be done by minimum distance, maximum likelihood, or the generalized method of moments.

Minimum distance is a straightforward application of the methods of Chapter 8 to the estimators presented in this chapter, as such methods apply to any asymptotically normal estimator.

Imposing restrictions on maximum likelihood is also straightforward. The likelihood is maximized subject to the imposed restrictions. One important example is explored in detail in the following section.

Generalized method of moments estimation of multivariate regression subject to restrictions will be explored in Section 13.18. This is a particularly simple and straightforward way to estimate restricted multivariate regression models and is our generally preferred approach.

11.11 Reduced Rank Regression

One context where systems estimation is important is when it is desired to impose or test restrictions across equations. Restricted systems are commonly estimated by maximum likelihood under normality. In this section we explore one important special case of restricted multivariate regression known as reduced rank regression. The model was originally proposed by Anderson (1951) and extended by Johansen (1995).

The unrestricted model is

$$\begin{aligned} Y &= \mathbf{B}'X + \mathbf{C}'Z + e \\ \mathbb{E}[ee' | X, Z] &= \Sigma \end{aligned} \tag{11.19}$$

where \mathbf{B} is $k \times m$, \mathbf{C} is $\ell \times m$, $Y \in \mathbb{R}^m$, $X \in \mathbb{R}^k$, and $Z \in \mathbb{R}^\ell$. We separate the regressors as X and Z because the coefficient matrix \mathbf{B} will be restricted while \mathbf{C} will be unrestricted.

The matrix \mathbf{B} is full rank if

$$\text{rank}(\mathbf{B}) = \min(k, m).$$

The reduced rank restriction is $\text{rank}(\mathbf{B}) = r < \min(k, m)$ for some known r .

The reduced rank restriction implies that we can write the coefficient matrix \mathbf{B} in the factored form $\mathbf{B} = \mathbf{G}\mathbf{A}'$ where \mathbf{A} is $m \times r$ and \mathbf{G} is $k \times r$. This representation is not unique as we can replace \mathbf{G} with $\mathbf{G}\mathbf{Q}$ and \mathbf{A} with $\mathbf{A}\mathbf{Q}^{-1'}$ for any invertible \mathbf{Q} and the same relation holds. Identification therefore requires a normalization of the coefficients. A conventional normalization is $\mathbf{G}'\mathbf{D}\mathbf{G} = \mathbf{I}_r$ for given \mathbf{D} .

Equivalently, the reduced rank restriction can be imposed by requiring that \mathbf{B} satisfy the restriction $\mathbf{B}\mathbf{A}_\perp = \mathbf{G}\mathbf{A}'\mathbf{A}_\perp = 0$ for some $m \times (m - r)$ coefficient matrix \mathbf{A}_\perp . Since \mathbf{G} is full rank this requires that $\mathbf{A}'\mathbf{A}_\perp = 0$, hence \mathbf{A}_\perp is the orthogonal complement of \mathbf{A} . Note that \mathbf{A}_\perp is not unique as it can be replaced by $\mathbf{A}_\perp\mathbf{Q}$ for any $(m - r) \times (m - r)$ invertible \mathbf{Q} . Thus if \mathbf{A}_\perp is to be estimated it requires a normalization.

We discuss methods for estimation of \mathbf{G} , \mathbf{A} , Σ , \mathbf{C} , and \mathbf{A}_\perp . The standard approach is maximum likelihood under the assumption that $e \sim N(0, \Sigma)$. The log-likelihood function for the sample is

$$\begin{aligned} \ell_n(\mathbf{G}, \mathbf{A}, \mathbf{C}, \Sigma) &= -\frac{nm}{2} \log(2\pi) - \frac{n}{2} \log(\det(\Sigma)) \\ &\quad - \frac{1}{2} \sum_{i=1}^n (Y_i - \mathbf{A}\mathbf{G}'X_i - \mathbf{C}'Z_i)' \Sigma^{-1} (Y_i - \mathbf{A}\mathbf{G}'X_i - \mathbf{C}'Z_i). \end{aligned}$$

Anderson (1951) derived the MLE by imposing the constraint $\mathbf{B}\mathbf{A}_\perp = 0$ via the method of Lagrange multipliers. This turns out to be algebraically cumbersome.

Johansen (1995) instead proposed the following straightforward concentration method. Treating \mathbf{G} as if it is known, maximize the log-likelihood with respect to the other parameters. Resubstituting these estimators we obtain the concentrated log-likelihood function with respect to \mathbf{G} . This can be maximized to find the MLE for \mathbf{G} . The other parameter estimators are then obtained by substitution. We now describe these steps in detail.

Given \mathbf{G} the likelihood is a normal multivariate regression in the variables $\mathbf{G}'\mathbf{X}$ and \mathbf{Z} , so the MLE for \mathbf{A} , \mathbf{C} and Σ are least squares. In particular, using the Frisch-Waugh-Lovell residual regression formula we can write the estimators for \mathbf{A} and Σ as

$$\hat{\mathbf{A}}(\mathbf{G}) = (\tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G}) (\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})^{-1}$$

and

$$\hat{\Sigma}(\mathbf{G}) = \frac{1}{n} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G} (\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})^{-1} \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}})$$

where $\tilde{\mathbf{Y}} = \mathbf{Y} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}$ and $\tilde{\mathbf{X}} = \mathbf{X} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}$.

Substituting these estimators into the log-likelihood function we obtain the concentrated likelihood function, which is a function of \mathbf{G} only.

$$\begin{aligned} \tilde{\ell}_n(\mathbf{G}) &= \ell_n(\mathbf{G}, \hat{\mathbf{A}}(\mathbf{G}), \hat{\mathbf{C}}(\mathbf{G}), \hat{\Sigma}(\mathbf{G})) \\ &= \frac{m}{2} (n \log(2\pi) - 1) - \frac{n}{2} \log \left[\det \left(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \mathbf{G} (\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})^{-1} \mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} \right) \right] \\ &= \frac{m}{2} (n \log(2\pi) - 1) - \frac{n}{2} \log \left(\det(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}}) \right) - \frac{n}{2} \log \left[\frac{\det \left(\mathbf{G}' (\tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}}) \mathbf{G} \right)}{\det(\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})} \right]. \end{aligned}$$

The third equality uses Theorem A.1.8. The MLE $\hat{\mathbf{G}}$ for \mathbf{G} is the maximizer of $\tilde{\ell}_n(\mathbf{G})$, or equivalently equals

$$\begin{aligned} \hat{\mathbf{G}} &= \underset{\mathbf{G}}{\operatorname{argmin}} \frac{\det \left(\mathbf{G}' (\tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}}) \mathbf{G} \right)}{\det(\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})} \\ &= \underset{\mathbf{G}}{\operatorname{argmax}} \frac{\det \left(\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}} \mathbf{G} \right)}{\det(\mathbf{G}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \mathbf{G})} \\ &= \{v_1, \dots, v_r\} \end{aligned} \tag{11.20}$$

which are the generalized eigenvectors of $\tilde{\mathbf{X}}' \tilde{\mathbf{Y}} (\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}})^{-1} \mathbf{Y}' \tilde{\mathbf{X}}$ with respect to $\tilde{\mathbf{X}}' \tilde{\mathbf{X}}$ corresponding to the r largest generalized eigenvalues. (Generalized eigenvalues and eigenvectors are discussed in Section A.14.) The estimator satisfies the normalization $\hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \hat{\mathbf{G}} = \mathbf{I}_r$. Letting v_j^* denote the eigenvectors of (11.20) we can also express $\hat{\mathbf{G}} = \{v_m^*, \dots, v_{m-r+1}^*\}$.

This is computationally straightforward. In MATLAB, for example, the generalized eigenvalues and eigenvectors of a matrix \mathbf{A} with respect to \mathbf{B} are found using the command `eig(A,B)`.

Given $\hat{\mathbf{G}}$, the MLE $\hat{\mathbf{A}}$, $\hat{\mathbf{C}}$, $\hat{\Sigma}$ are found by least squares regression of \mathbf{Y} on $\hat{\mathbf{G}}'\mathbf{X}$ and \mathbf{Z} . In particular, $\hat{\mathbf{A}} = \hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{Y}}$ because $\hat{\mathbf{G}}' \tilde{\mathbf{X}}' \tilde{\mathbf{X}} \hat{\mathbf{G}} = \mathbf{I}_r$.

We now discuss the estimator $\hat{\mathbf{A}}_{\perp}$ of \mathbf{A}_{\perp} . It turns out that

$$\begin{aligned}\hat{\mathbf{A}}_{\perp} &= \underset{\mathbf{A}}{\operatorname{argmax}} \frac{\det \left(\mathbf{A}' \left(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}' \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} \right) \mathbf{A} \right)}{\det \left(\mathbf{A}' \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \mathbf{A} \right)} \\ &= \{w_1, \dots, w_{m-r}\}\end{aligned}\tag{11.21}$$

the eigenvectors of $\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}' \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}}$ with respect to $\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}}$ associated with the largest $m - r$ eigenvalues.

By the dual eigenvalue relation (Theorem A.5), equations (11.20) and (11.21) have the same non-zero eigenvalues λ_j and the associated eigenvectors v_j^* and w_j satisfy the relationship

$$w_j = \lambda_j^{-1/2} \left(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \right)^{-1} \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} v_j^*.$$

Letting $\Lambda = \operatorname{diag}\{\lambda_m, \dots, \lambda_{m-r+1}\}$ this implies

$$\{w_m, \dots, w_{m-r+1}\} = \left(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \right)^{-1} \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \{v_m^*, \dots, v_{m-r+1}^*\} \Lambda = \left(\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \right)^{-1} \hat{\mathbf{A}} \Lambda.$$

The second equality holds because $\hat{\mathbf{G}} = \{v_m^*, \dots, v_{m-r+1}^*\}$ and $\hat{\mathbf{A}} = \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \hat{\mathbf{G}}$. Since the eigenvectors w_j satisfy the orthogonality property $w_j' \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} w_{\ell} = 0$ for $j \neq \ell$, it follows that

$$0 = \hat{\mathbf{A}}'_{\perp} \tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} \{w_m, \dots, w_{m-r+1}\} = \hat{\mathbf{A}}'_{\perp} \hat{\mathbf{A}} \Lambda.$$

Since $\Lambda > 0$ we conclude that $\hat{\mathbf{A}}'_{\perp} \hat{\mathbf{A}} = 0$ as desired.

The solution $\hat{\mathbf{A}}_{\perp}$ in (11.21) can be represented several ways. One which is computationally convenient is to observe that

$$\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}} - \tilde{\mathbf{Y}}' \tilde{\mathbf{X}} \left(\tilde{\mathbf{X}}' \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{Y}} = \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} = \tilde{\mathbf{E}}' \tilde{\mathbf{E}}$$

where $\mathbf{M}_{\mathbf{X}, \mathbf{Z}} = \mathbf{I}_n - (\mathbf{X}, \mathbf{Z}) \left((\mathbf{X}, \mathbf{Z})' (\mathbf{X}, \mathbf{Z}) \right)^{-1} (\mathbf{X}, \mathbf{Z})'$ and $\tilde{\mathbf{E}} = \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y}$ is the residual matrix from the unrestricted multivariate least squares regression of \mathbf{Y} on \mathbf{X} and \mathbf{Z} . The first equality follows by the Frisch-Waugh-Lovell theorem. This shows that $\hat{\mathbf{A}}_{\perp}$ are the generalized eigenvectors of $\tilde{\mathbf{E}}' \tilde{\mathbf{E}}$ with respect to $\tilde{\mathbf{Y}}' \tilde{\mathbf{Y}}$ corresponding to the $m - r$ largest eigenvalues. In MATLAB, for example, these can be computed using the `eig(A, B)` command.

Another representation is to write $\mathbf{M}_{\mathbf{Z}} = \mathbf{I}_n - \mathbf{Z} (\mathbf{Z}' \mathbf{Z})^{-1} \mathbf{Z}'$ so that

$$\hat{\mathbf{A}}_{\perp} = \underset{\mathbf{A}}{\operatorname{argmax}} \frac{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} \mathbf{A})}{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{Z}} \mathbf{Y} \mathbf{A})} = \underset{\mathbf{A}}{\operatorname{argmin}} \frac{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{Z}} \mathbf{Y} \mathbf{A})}{\det(\mathbf{A}' \mathbf{Y}' \mathbf{M}_{\mathbf{X}, \mathbf{Z}} \mathbf{Y} \mathbf{A})}.$$

We summarize our findings.

Theorem 11.7 The MLE for the reduced rank model (11.19) under $e \sim N(0, \Sigma)$ is given as follows. Let \tilde{Y} and \tilde{X} be the residual matrices from multivariate regression of Y and X on Z , respectively. Then $\hat{G}_{\text{mle}} = \{v_1, \dots, v_r\}$, the generalized eigenvectors of $\tilde{X}'\tilde{Y}(\tilde{Y}'\tilde{Y})^{-1}\tilde{Y}'\tilde{X}$ with respect to $\tilde{X}'\tilde{X}$ corresponding to the r largest eigenvalues $\hat{\lambda}_j$. \hat{A}_{mle} , \hat{C}_{mle} and $\hat{\Sigma}_{\text{mle}}$ are obtained by the least squares regression

$$Y_i = \hat{A}_{\text{mle}} \hat{G}_{\text{mle}}' X_i + \hat{C}_{\text{mle}}' Z_i + \hat{e}_i$$

$$\hat{\Sigma}_{\text{mle}} = \frac{1}{n} \sum_{i=1}^n \hat{e}_i \hat{e}_i'.$$

Let \tilde{E} be the residual matrix from a multivariate regression of Y on X and Z . Then \hat{A}_{\perp} equals the generalized eigenvectors of $\tilde{E}'\tilde{E}$ with respect to $\tilde{Y}'\tilde{Y}$ corresponding to the $m - r$ smallest eigenvalues. The maximized likelihood equals

$$\ell_n = \frac{m}{2} (n \log(2\pi) - 1) - \frac{n}{2} \log(\det(\tilde{Y}'\tilde{Y})) - \frac{n}{2} \sum_{j=1}^r \log(1 - \hat{\lambda}_j).$$

An R package for reduced rank regression is “RRR”. I am unaware of a Stata command.

11.12 Principal Component Analysis

In Section 4.21 we described the Duflo, Dupas, and Kremer (2011) dataset which is a sample of Kenyan first grade test scores. Following the authors we focused on the variable *totalscore* which is each student's composite test score. If you examine the data file you will find other pieces of information about the students' performance, including each student's score on separate sections of the test, with the labels *wordscore* (word recognition), *sentscore* (sentence recognition), *letterscore* (letter recognition), *spellscore* (spelling), *additions_score* (addition), *subtractions_score* (subtraction), *multiplications_score* (multiplication). The “total” score sums the scores from the individual sections. Perhaps there is more information in the section scores. How can we learn about this from the data?

Principal component analysis (PCA) addresses this issue by ordering linear combinations by their contribution to variance.

Definition 11.1 Let X be a $k \times 1$ random vector.

The **first principal component** is $U_1 = h'_1 X$ where h_1 satisfies

$$h_1 = \operatorname{argmax}_{h'h=1} \operatorname{var}[h'X].$$

The **second principal component** is $U_2 = h'_2 X$ where

$$h_2 = \operatorname{argmax}_{h'h=1, h'h_1=0} \operatorname{var}[h'X].$$

In general, the j^{th} **principal component** is $U_j = h'_j X$ where

$$h_j = \operatorname{argmax}_{h'h=1, h'h_1=0, \dots, h'h_{j-1}=0} \operatorname{var}[h'X].$$

The principal components of X are linear combinations $h'X$ ranked by contribution to variance. By the properties of quadratic forms (Section A.15) the weight vectors h_j are the eigenvectors of $\Sigma = \operatorname{var}[X]$.

Theorem 11.8 The principal components of X are $U_j = h'_j X$, where h_j is the eigenvector of Σ associated with the j^{th} ordered eigenvalue λ_j of Σ .

Another way to see the PCA construction is as follows. Since Σ is symmetric the spectral decomposition (Theorem A.3) states that $\Sigma = \mathbf{H}\mathbf{D}\mathbf{H}'$ where $\mathbf{H} = [h_1, \dots, h_k]$ and $\mathbf{D} = \operatorname{diag}(d_1, \dots, d_k)$ are the eigenvectors and eigenvalues of Σ . Since Σ is positive semi-definite the eigenvalues are real, non-negative, and ordered $d_1 \geq d_2 \geq \dots \geq d_k$. Let $U = (U_1, \dots, U_k)$ be the principal components of X . By Theorem 11.8, $U = \mathbf{H}'X$. The covariance matrix of U is

$$\operatorname{var}[U] = \operatorname{var}[\mathbf{H}'X] = \mathbf{H}'\Sigma\mathbf{H} = \mathbf{D}$$

which is diagonal. This shows that $\operatorname{var}[U_j] = d_j$ and the principal components are mutually uncorrelated. The relative variance contribution of the j^{th} principal component is $d_j / \operatorname{tr}(\Sigma)$.

Principal components are sensitive to the scaling of X . Consequently, it is recommended to first scale each element of X to have mean zero and unit variance. In this case Σ is a correlation matrix.

The sample principal components are obtained by replacing the unknowns by sample estimators. Let $\hat{\Sigma}$ be the sample covariance or correlation matrix and $\hat{h}_1, \hat{h}_2, \dots, \hat{h}_k$ its ordered eigenvectors. The sample principal components are $\hat{h}'_j X_i$.

To illustrate we use the Duflo, Dupas, and Kremer (2011) dataset. In Table 11.1 we display the seven eigenvalues of the sample correlation matrix for the seven test scores described above. The seven eigenvalues sum to seven because we have applied PCA to the correlation matrix. The first eigenvalue is 4.0, implying that the first principal component explains 57% of the variance of the seven test scores. The second eigenvalue is 1.0, implying that the second principal component explains 15% of the variance. Together the first two components explain 72% of the variance of the seven test scores.

In Table 11.2 we display the weight vectors (eigenvectors) for the first two principal components. The weights for the first component are all positive and similar in magnitude. This means that the first

Table 11.1: Eigenvalue Decomposition of Sample Correlation Matrix

	Eigenvalue	Proportion
1	4.02	0.57
2	1.04	0.15
3	0.57	0.08
4	0.52	0.08
5	0.37	0.05
6	0.29	0.04
7	0.19	0.03

Table 11.2: Principal Component Weight Vectors

	First	Second
words	0.41	-0.32
sentences	0.32	-0.49
letters	0.40	-0.13
spelling	0.43	-0.28
addition	0.38	0.41
subtraction	0.35	0.52
multiplication	0.33	0.36

principal component is similar to a simple average of the seven test scores. This is quite fascinating. This is consistent with our intuition that a simple average (e.g. the variable *totalscore*) captures most of the information contained in the seven test scores. The weights for the second component have a different pattern. The four literacy scores receive negative weight and the three math scores receive positive weight with similar magnitudes. This means that the second principal component is similar to the difference between a student's math and verbal test scores. Taken together, the information in the first two principal components is equivalent to "average verbal" and "average math" test scores. What this shows is that 57% of the variation in the seven section test scores can be explained by a simple average (e.g. *totalscore*), and 72% can be explained by averages for the verbal and math halves of the test.

In Stata, principal components analysis can be implemented with the `pca` command. In R use `prcomp` or `princomp`. All three can be applied to either covariance matrices (unscaled data) or correlation matrices (normalized data) but they have different default settings. The Stata `pca` command by default normalizes the observations. The R commands by default do not normalize the observations.

11.13 Factor Models

Closely related to principal components are factor models. These are statistical models which decompose random vectors into common factors and idiosyncratic errors. Factor models are popular throughout the social sciences. Consequently a variety of estimation methods have been developed. In the next few sections we focus on methods which are popular among economists.

Let $X = (X_1, \dots, X_k)'$ be a $k \times 1$ random vector (for example the seven test scores described in the previous section). Assume that the elements of X are scaled to have mean zero and unit variance.

A **single factor model** for X is

$$X = \lambda F + u \quad (11.22)$$

where $\lambda \in \mathbb{R}^k$ are **factor loadings**, $F \in \mathbb{R}$ is a **common factor**, and $u \in \mathbb{R}^k$ is a random error. The factor F is individual-specific while the coefficient λ is common across individuals. The model (11.22) specifies that correlation between the elements of X is due to the common factor F . In the student test score example it is intuitive to think of F as a student's scholastic "aptitude"; in this case the vector λ describes how scholastic aptitude affects the seven subject scores.

A multiple factor model has $r < k$ factors. We write the model as

$$X = \Lambda F + u \quad (11.23)$$

where Λ is a $k \times r$ matrix of factor loadings and $F = (F_1, \dots, F_r)'$ is an $r \times 1$ vector of factors. In the student test score example possible factors could be "math aptitude", "language skills", "social skills", "artistic ability", "creativity", etc. The factor loading matrix Λ indicates the effect of each factor on each test score. The number of factors r is taken as known. We discuss selection of r later.

The error vector u is assumed to be mean zero, uncorrelated with F , and (under correct specification) to have mutually uncorrelated elements. We write its covariance matrix as $\Psi = \mathbb{E}[uu']$. The factor vector F can either be treated as random or as a regressor. In this section we treat F as random; in the next we treat F as regressors. The random factors F are assumed mean zero and are normalized so that $\mathbb{E}[FF'] = I_r$.

The assumptions imply that the correlation matrix $\Sigma = \mathbb{E}[XX']$ equals

$$\Sigma = \Lambda\Lambda' + \Psi. \quad (11.24)$$

The factor analysis literature often describes $\Lambda\Lambda'$ as the **communality** and the idiosyncratic error matrix Ψ as the **uniqueness**. The former is the portion of the variance which is explained by the factor model and the latter is the unexplained portion of the variance.

The model is often¹ estimated by **maximum likelihood**. Under joint normality of (F, u) the distribution of X is $N(0, \Lambda\Lambda' + \Psi)$. The parameters are Λ and $\Psi = \text{diag}(\psi_1, \dots, \psi_k)$. The log-likelihood function of a random sample (X_1, \dots, X_n) is

$$\ell_n(\Lambda, \Psi) = -\frac{nk}{2} \log(2\pi) - \frac{n}{2} \log \det(\Lambda\Lambda' + \Psi) - \frac{n}{2} \text{tr}\left((\Lambda\Lambda' + \Psi)^{-1} \hat{\Sigma}\right). \quad (11.25)$$

The MLE $(\hat{\Lambda}, \hat{\Psi})$ maximizes $\ell_n(\Lambda, \Psi)$. There is not an algebraic solution so the estimator is found using numerical methods. Fortunately, computational algorithms are available in standard packages. A detailed description and analysis can be found in Anderson (2003, Chapter 14).

The form of the log-likelihood is intriguing. Notice that the log-likelihood is only a function of the observations through its correlation matrix $\hat{\Sigma}$, and only a function of the parameters through the population correlation matrix $\Lambda\Lambda' + \Psi$. The final term in (11.25) is a measure of the match between $\hat{\Sigma}$ and $\Lambda\Lambda' + \Psi$. Together, we see that the Gaussian log-likelihood is essentially a measure of the fit of the model and sample correlation matrices. It is therefore not reliant on the normality assumption.

It is often of interest to estimate the factors F_i . Given Λ the equation $X_i = \Lambda F_i + u_i$ can be viewed as a regression with coefficient F_i . Its least squares estimator is $\hat{F}_i = (\Lambda'\Lambda)^{-1} \Lambda' X_i$. The GLS estimator (taking into account the covariance matrix of u_i) is $\tilde{F}_i = (\Lambda'\Psi^{-1}\Lambda)^{-1} \Lambda'\Psi^{-1} X_i$. This motivates the **Bartlett scoring** estimator

$$\tilde{F}_i = (\hat{\Lambda}'\hat{\Psi}^{-1}\hat{\Lambda})^{-1} \hat{\Lambda}'\hat{\Psi}^{-1} X_i.$$

The idealized version satisfies

$$\hat{F}_i = (\Lambda'\Psi^{-1}\Lambda)^{-1} \Lambda'\Psi^{-1} (\Lambda F_i + u_i) = F_i + (\Lambda'\Psi^{-1}\Lambda)^{-1} \Lambda'\Psi^{-1} u_i$$

¹There are other estimators used in applied factor analysis. However there is little reason to consider estimators beyond the MLE of this section and the principal components estimator of the next section.

which is unbiased for F_i and has variance $(\Lambda' \Psi^{-1} \Lambda)^{-1}$. Thus the Bartlett scoring estimator is typically described as “unbiased” though this is actually a property of its idealized version \hat{F}_i .

A second estimator for the factors can be constructed from the multivariate linear projection of F on X . This is $F = \mathbf{A}X + \xi$ where the coefficient matrix \mathbf{A} is $r \times k$. The coefficient matrix equals

$$\mathbf{A} = \mathbb{E}[FX'] \mathbb{E}[XX']^{-1} = \Lambda' \Sigma^{-1},$$

the second equation using $\mathbb{E}[FX'] = \mathbb{E}[F(\Lambda F + u)'] = \mathbb{E}[FF'] \Lambda' + \mathbb{E}[Fu'] = \Lambda'$. The predicted value of F_i is $F_i^* = \mathbf{A}X_i = \Lambda' \Sigma^{-1} X_i$. This motivates the **regression scoring** estimator

$$\bar{F}_i = \hat{\Lambda}' \hat{\Sigma}^{-1} X_i.$$

The idealized version F_i^* has conditional expectation $\Lambda' \Sigma^{-1} \Lambda F_i$ and is thus biased for F_i . Hence the regression scoring estimator \bar{F}_i is often described as “biased”. Some algebraic manipulations reveal that F_i^* has MSE $\mathbf{I}_r - \Lambda' (\Lambda' \Lambda + \Psi)^{-1} \Lambda$ which is smaller (in a positive definite sense) than the MSE of the idealized Bartlett estimator \hat{F}_i .

Which estimator is preferred, Bartlett or regression scoring? The differences diminish when k is large so the choice is most relevant for small to moderate k . The regression scoring estimator has lower approximate MSE, meaning that it is a more precise estimator. Thus based on estimation precision this is our recommended choice.

The factor loadings Λ and factors F are not separately identified. To see this, notice that if you replace (Λ, F) with $\Lambda^* = \Lambda \mathbf{G}$ and $F^* = \mathbf{G}' F$ where \mathbf{G} is $r \times r$ and orthonormal then the regression model is identical. Such replacements are called “rotations” in the factor analysis literature. Any orthogonal rotation of the factor loadings is an equally valid representation. The default MLE outputs are one specific rotation; others can be obtained by a variety of algorithms (which we do not review here). Consequently it is unwise to attribute meaning to the individual factor loading estimates.

Another important and tricky issue is selection of the number of factors r . There is no clear guideline. One approach is to examine the principal component decomposition, look for a division between the “large” and “small eigenvalues, and set r to equal to the number of “large” eigenvalues. Another approach is based on testing. As a by-product of the MLE (and standard package implementations) we obtain the LR test for the null hypothesis of r factors against the alternative hypothesis of k factors. If the LR test rejects (has a small p-value) this is evidence that the given r may be too small.

In Stata, the MLE $(\hat{\Lambda}, \hat{\Psi})$ can be calculated with the factor, ml factors(r) command. The factor estimates \tilde{F}_i and \bar{F}_i can be calculated by the predict command with either the barlett or regression option, respectively. In R, the command `factanal(X, factors=r, rotation="none")` calculates the MLE $(\hat{\Lambda}, \hat{\Psi})$ and also calculates the factor estimates \tilde{F}_i and/or \bar{F}_i using the scores option.

11.14 Approximate Factor Models

The MLE of the previous section is a good choice for factor estimation when the number of variables k is small and the factor model is believed to be correctly specified. In many economic applications of factor analysis, however, the number of variables k is large. In such contexts the MLE can be computationally costly and/or unstable. Furthermore it is typically not credible to believe that the model is correctly specified; rather it is more reasonable to view the factor model as a useful approximation. In this section we explore an approach known as the approximate factor model with estimation by principal components. The estimation method is justified by an asymptotic framework where the number of variables $k \rightarrow \infty$.

The **approximate factor model** was introduced by Chamberlain and Rothschild (1983). It is the same as (11.23) but relaxes the assumption on the idiosyncratic error u so that the covariance matrix $\Psi = \mathbb{E}[uu']$ is left unrestricted. In this context the Gaussian MLE of the previous section is misspecified.

Chamberlain and Rothschild (and the literature which followed) proposed estimation by least squares. The idea is to treat the factors as unknown regressors and simultaneously estimate the factors F_i and factor loadings Λ . We first describe the estimation method.

Let (X_1, \dots, X_n) be a sample centered at sample means. The least squares criterion is

$$\frac{1}{n} \sum_{i=1}^n (X_i - \Lambda F_i)' (X_i - \Lambda F_i).$$

Let $(\hat{\Lambda}, \hat{F}_1, \dots, \hat{F}_n)$ be the joint minimizers. As Λ and F_i are not separately identified a normalization is needed. For compatibility with the notation of the previous section we use $n^{-1} \sum_{i=1}^n \hat{F}_i \hat{F}_i' = I_r$.

We use a concentration argument to find the solution. As described in the previous section, each observation satisfies the multivariate equation $X_i = \Lambda F_i + u_i$. For fixed Λ this is a set of k equations with r unknowns F_i . The least squares solution is $\hat{F}_i(\Lambda) = (\Lambda' \Lambda)^{-1} \Lambda' X_i$. Substituting this expression into the least squares criterion the concentrated least squares criterion for Λ is

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (X_i - \Lambda \hat{F}_i(\Lambda))' (X_i - \Lambda \hat{F}_i(\Lambda)) &= \frac{1}{n} \sum_{i=1}^n \left(X_i - \Lambda (\Lambda' \Lambda)^{-1} \Lambda' X_i \right)' \left(X_i - \Lambda (\Lambda' \Lambda)^{-1} \Lambda' X_i \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left(X_i' X_i - X_i' \Lambda (\Lambda' \Lambda)^{-1} \Lambda' X_i \right) \\ &= \text{tr}[\hat{\Sigma}] - \text{tr}[(\Lambda' \Lambda)^{-1} \Lambda' \hat{\Sigma} \Lambda] \end{aligned}$$

where $\hat{\Sigma} = n^{-1} \sum_{i=1}^n X_i X_i'$ is the sample covariance matrix. The least squares estimator $\hat{\Lambda}$ minimizes this criterion. Let \hat{D} and \hat{H} be first r eigenvalues and eigenvectors of $\hat{\Sigma}$. Using the normalization $\Lambda' \Lambda = I_r$, from the extrema results of Section A.15 the minimizer of the least squares criterion is $\hat{\Lambda} = \hat{H}$. More broadly any rotation of \hat{H} is valid. Consider $\hat{\Lambda} = \hat{H} \hat{D}^{1/2}$. Recall the expression for the factors $\hat{F}_i(\Lambda) = (\Lambda' \Lambda)^{-1} \Lambda' X_i$. We find that the estimated factors are

$$\hat{F}_i = \left(\hat{D}^{1/2} \hat{H}' \hat{H} \hat{D}^{1/2} \right)^{-1} \hat{D}^{1/2} \hat{H}' X_i = \hat{D}^{-1/2} \hat{H}' X_i.$$

We calculate that

$$n^{-1} \sum_{i=1}^n \hat{F}_i \hat{F}_i' = \hat{D}^{-1/2} \hat{H}' \hat{\Sigma} \hat{H} \hat{D}^{-1/2} = \hat{D}^{-1/2} \hat{D} \hat{D}^{-1/2} = I_r$$

which is the desired normalization. This shows that the rotation $\hat{\Lambda} = \hat{H} \hat{D}^{1/2}$ produces factor estimates satisfying this normalization.

We have proven the following result.

Theorem 11.9 The least squares estimator of the factor model (11.23) under the normalization $n^{-1} \sum_{i=1}^n \hat{F}_i \hat{F}_i' = I_r$ has the following solution:

1. Let $\hat{D} = \text{diag}[\hat{d}_1, \dots, \hat{d}_r]$ and $\hat{H} = [\hat{h}_1, \dots, \hat{h}_r]$ be the first r eigenvalues and eigenvectors of the sample covariance matrix $\hat{\Sigma}$.
2. $\hat{\Lambda} = \hat{H} \hat{D}^{1/2}$.
3. $\hat{F}_i = \hat{D}^{-1/2} \hat{H}' X_i$.

Theorem 11.9 shows that the least squares estimator is based on an eigenvalue decomposition of the covariance matrix. This is computationally stable even in high dimensions.

The factor estimates are the principal components scaled by the eigenvalues of $\hat{\Sigma}$. Specifically, the j^{th} factor estimate is $\hat{F}_{ji} = \hat{d}_j^{-1/2} \hat{h}_j' X_i$. Consequently many authors call this estimator the “principal-component method”.

Unfortunately, $\hat{\Lambda}$ is inconsistent for Λ if k is fixed, as we now show. By the WLLN and CMT, $\hat{\Sigma} \xrightarrow{p} \Sigma$ and $\hat{H} \xrightarrow{p} H$, the first r eigenvectors of Σ . When Ψ is diagonal, the eigenvectors of $\Sigma = \Lambda\Lambda' + \Psi$ do not lie in the range space of Λ except in the special case $\Psi = \sigma^2 I_k$. Consequently the estimator $\hat{\Lambda}$ is inconsistent.

This inconsistency should not be viewed as surprising. The sample has a total of nk observations and the model has a total of $nr + kr - r(r+1)/2$ parameters. Since the number of estimated parameters is proportional to sample size we should not expect estimator consistency.

As first recognized by Chamberlain and Rothschild, this deficiency diminishes as k increases. Specifically, assume that $k \rightarrow \infty$ as $n \rightarrow \infty$. One implication is that the number of observations nk increase at a rate faster than n , while the number of parameters increase at a rate proportional to n . Another implication is that as k increases there is increasing information about the factors.

To make this precise we add the following assumption. Let $\lambda_{\min}(A)$ and $\lambda_{\max}(A)$ denote the smallest and largest eigenvalues of a positive semi-definite matrix A .

Assumption 11.1 As $k \rightarrow \infty$

1. $\lambda_{\max}(\Psi) \leq B < \infty$.
2. $\lambda_{\min}(\Lambda'\Lambda) \rightarrow \infty$ as $k \rightarrow \infty$.

Assumption 11.1.1 bounds the covariance matrix of the idiosyncratic errors. When $\Psi = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ this is the same as bounding the individual variances. Effectively Assumption 11.1.1 means that while the elements of u can be correlated they cannot have a correlation structure similar to that of a factor model. Assumption 11.1.2 requires the factor loading matrix to increase in magnitude as the number of variables increases. This is a fairly mild requirement. When the factor loadings are of similar magnitude across variables, $\lambda_{\min}(\Lambda'\Lambda) \sim k \rightarrow \infty$. Conceptually, Assumption 11.1.2 requires additional variables to add information about the unobserved factors.

Assumption 11.1 implies that in the covariance matrix factorization $\Sigma = \Lambda\Lambda' + \Psi$ the component $\Lambda\Lambda'$ dominates as k increases. This means that for large k the first r eigenvectors of Σ are equivalent to those of $\Lambda\Lambda'$, which are in the range space of Λ . This observation led Chamberlain and Rothschild (1983) to deduce that the principal components estimator is an asymptotic (large k) analog estimator for the factor loadings and factors. Bai (2003) demonstrated that the estimator is consistent as $n, k \rightarrow \infty$ jointly. The conditions and proofs are technical so are not reviewed here.

Now consider the estimated factors

$$\hat{F}_i = D^{-1/2} H' X_i = D^{-1} \Lambda' X_i$$

where for simplicity we ignore estimation error. Since $X_i = \Lambda F_i + u_i$ and $\Lambda'\Lambda = D$ we can write this as

$$\hat{F}_i = F_i + D^{-1} \Lambda' u_i.$$

This shows that \hat{F}_i is an unbiased estimator for F_i and has variance $\text{var}[\hat{F}_i] = \mathbf{D}^{-1}\Lambda'\Psi\Lambda\mathbf{D}^{-1}$. Under Assumption 11.1, $\|\text{var}[\hat{F}_i]\| \leq B/\lambda_{\min}(\Lambda'\Lambda) \rightarrow 0$. Thus \hat{F}_i is consistent for F_i as $k \rightarrow \infty$. Bai (2003) shows that this extends to the feasible estimator as $n, k \rightarrow \infty$.

In Stata, the least squares estimator $\hat{\Lambda}$ and factors \hat{F}_i can be calculated with the `factor, pcf factors(r)` command followed by `predict`. In R a feasible estimation approach is to calculate the factors by eigenvalue decomposition.

11.15 Factor Models with Additional Regressors

Consider the model

$$X = \Lambda F + \mathbf{B}Z + e$$

where X and e are $k \times 1$, Λ is $k \times r$, F is $r \times 1$, \mathbf{B} is $k \times \ell$, and Z is $\ell \times 1$.

The coefficients Λ and \mathbf{B} can be estimated by a combination of factor regression (either MLE or principal components) and least squares. The key is the following two observations:

1. Given \mathbf{B} , the coefficient Λ can be estimated by factor regression applied to $X - \mathbf{B}Z$.
2. Given the factors F , the coefficients Λ and \mathbf{B} can be estimated by multivariate least squares of X on F and Z .

Estimation iterates between these two steps. Start with a preliminary estimator of \mathbf{B} obtained by multivariate least squares of X on Z . Then apply the above two steps and iterate under convergence.

11.16 Factor-Augmented Regression

In the previous sections we considered factor models which decompose a set of variables into common factors and idiosyncratic errors. In this section we consider factor-augmented regression, which uses such common factors as regressors for dimension reduction.

Suppose we have the variables (Y, Z, X) where $Y \in \mathbb{R}$, $Z \in \mathbb{R}^\ell$, and $X \in \mathbb{R}^k$. In practice, k may be large and the elements of X may be highly correlated. The **factor-augmented regression** model is

$$\begin{aligned} Y &= F'\beta + Z'\gamma + e \\ X &= \Lambda F + u \\ \mathbb{E}[Fe] &= 0 \\ \mathbb{E}[Ze] &= 0 \\ \mathbb{E}[Fu'] &= 0 \\ \mathbb{E}[ue] &= 0, \end{aligned}$$

The random variables are $e \in \mathbb{R}$, $F \in \mathbb{R}^r$, and $u \in \mathbb{R}^k$. The regression coefficients are $\beta \in \mathbb{R}^k$ and $\gamma \in \mathbb{R}^\ell$. The matrix Λ are the factor loadings.

This model specifies that the influence of X on Y is through the common factors F . The idea is that the variation in the regressors is mostly captured by the variation in the factors, so the influence of the regressors can be captured through these factors. This can be viewed as a dimension-reduction technique as we have reduced the k -dimensional X to the r -dimensional F . Interest typically focuses on the regressors Z and its coefficients γ . The factors F are included in the regression as “controls” and its coefficient β is less typically of interest. Since it is difficult to interpret the factors F only their range space is identified it is generally prudent to avoid interpreting the coefficients β .

The model is typically estimated in multiple steps. First, the factor loadings Λ and factors F_i are estimated by factor regression. In the case of principal-components estimation the factor estimates are the scaled² principal components $\hat{F}_i = \hat{\mathbf{D}}^{-1} \hat{\Lambda}' X_i$. Second, Y is regressed on the estimated factors and the other regressors to obtain the estimator of β and γ . This second-step estimator equals (for simplicity assume there is no Z)

$$\begin{aligned} \hat{\beta} &= \left(\sum_{i=1}^n \hat{F}_i \hat{F}_i' \right)^{-1} \left(\sum_{i=1}^n \hat{F}_i Y_i \right) \\ &= \left(\hat{\mathbf{D}}^{-1} \hat{\Lambda}' \frac{1}{n} \sum_{i=1}^n X_i X_i' \hat{\Lambda} \hat{\mathbf{D}}^{-1} \right)^{-1} \left(\hat{\mathbf{D}}^{-1} \hat{\Lambda}' \frac{1}{n} \sum_{i=1}^n X_i Y_i \right). \end{aligned}$$

Now let's investigate its asymptotic behavior. As $n \rightarrow \infty$, $\hat{\Lambda} \xrightarrow{p} \Lambda$ and $\hat{\mathbf{D}} \xrightarrow{p} \mathbf{D}$ so

$$\hat{\beta} \xrightarrow{p} \beta^* = (\mathbf{D}^{-1} \Lambda' \mathbb{E}[X X'] \Lambda \mathbf{D}^{-1})^{-1} (\mathbf{D}^{-1} \Lambda' \mathbb{E}[X Y]). \quad (11.26)$$

Recall $\mathbb{E}[X X'] = \Lambda \Lambda' + \Psi$ and $\Lambda' \Lambda = \mathbf{D}$. We calculate that

$$\mathbb{E}[X Y] = \mathbb{E}[(\Lambda F + u)(F' \beta + e)] = \Lambda \beta.$$

We find that the right-hand-side of (11.26) equals

$$\beta^* = (\mathbf{D}^{-1} \Lambda' (\Lambda \Lambda' + \Psi) \Lambda \mathbf{D}^{-1})^{-1} (\mathbf{D}^{-1} \Lambda' \Lambda \beta) = (\mathbf{I}_r + \mathbf{D}^{-1} \Lambda' \Psi \Lambda \mathbf{D}^{-1})^{-1} \beta$$

which does not equal β . Thus $\hat{\beta}$ has a probability limit but is inconsistent for β as $n \rightarrow \infty$.

This deficiency diminishes as $k \rightarrow \infty$. Indeed,

$$\|\mathbf{D}^{-1} \Lambda' \Psi \Lambda \mathbf{D}^{-1}\| \leq B \|\mathbf{D}^{-1}\| \rightarrow 0$$

as $k \rightarrow \infty$. This implies $\beta^* \rightarrow \beta$. Hence, if we take the sequential asymptotic limit $n \rightarrow \infty$ followed by $k \rightarrow \infty$, we find $\hat{\beta} \xrightarrow{p} \beta$. This implies that the estimator is consistent. Bai (2003) demonstrated consistency under the more rigorous but technically challenging setting where $n, k \rightarrow \infty$ jointly. The implication of this result is that factor augmented regression is consistent if both the sample size and dimension of X are large.

For asymptotic normality of $\hat{\beta}$ it turns out that we need to strengthen Assumption 11.1.2. The relevant condition is $n^{-1/2} \lambda_{\min}(\Lambda' \Lambda) \rightarrow \infty$. This is similar to the condition that $k^2/n \rightarrow \infty$. This is technical but can be interpreted as meaning that k is large relative to \sqrt{n} . Intuitively, this requires that the dimension of X is larger than sample size n .

In Stata, estimation takes the following steps. First, the `factor` command is used to estimate the factor model. Either MLE or principal components estimation can be used. Second, the `predict` command is used to estimate the factors, either by Barlett or regression scoring. Third, the factors are treated as regressors in an estimated regression.

11.17 Multivariate Normal*

Some interesting sampling results hold for matrix-valued normal variates. Let \mathbf{Y} be an $n \times m$ matrix whose rows are independent and distributed $N(\mu, \Sigma)$. We say that \mathbf{Y} is **multivariate matrix normal**, and

²The unscaled principal components can equivalently be used if the coefficients $\hat{\beta}$ are not reported. The coefficient estimates $\hat{\gamma}$ are unaffected by the choice of factor scaling.

write $\mathbf{Y} \sim N(\bar{\boldsymbol{\mu}}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$, where $\bar{\boldsymbol{\mu}}$ is $n \times m$ with each row equal to $\boldsymbol{\mu}'$. The notation is due to the fact that $\text{vec}((\mathbf{Y} - \boldsymbol{\mu})') \sim N(0, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$.

Definition 11.2 If $n \times m$ $\mathbf{Y} \sim N(\bar{\boldsymbol{\mu}}, \mathbf{I}_n \otimes \boldsymbol{\Sigma})$ then $\mathbf{W} = \mathbf{Y}'\mathbf{Y}$ is distributed **Wishart** with n degrees of freedom and covariance matrix $\boldsymbol{\Sigma}$, and is written as $\mathbf{W} \sim W_m(n, \boldsymbol{\Sigma})$.

The Wishart is a multivariate generalization of the chi-square. If $\mathbf{W} \sim W_1(n, \sigma^2)$ then $\mathbf{W} \sim \sigma^2 \chi_n^2$.

The Wishart arises as the exact distribution of a sample covariance matrix in the normal sampling model. The bias-corrected estimator of $\boldsymbol{\Sigma}$ is

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{Y}_i - \bar{\mathbf{Y}})(\mathbf{Y}_i - \bar{\mathbf{Y}})'$$

Theorem 11.10 If $\mathbf{Y}_i \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ are independent then $\hat{\boldsymbol{\Sigma}} \sim W_m(n-1, \frac{1}{n-1}\boldsymbol{\Sigma})$.

The following manipulation is useful.

Theorem 11.11 If $\mathbf{W} \sim W_m(n, \boldsymbol{\Sigma})$ then for $m \times 1$ $\boldsymbol{\alpha}$, $(\boldsymbol{\alpha}'\mathbf{W}^{-1}\boldsymbol{\alpha})^{-1} \sim \frac{\chi_{n-m+1}^2}{\boldsymbol{\alpha}'\boldsymbol{\Sigma}^{-1}\boldsymbol{\alpha}}$.

To prove this, note that without loss of generality we can take $\boldsymbol{\Sigma} = \mathbf{I}_m$ and $\boldsymbol{\alpha}'\boldsymbol{\alpha} = 1$. Let \mathbf{H} be $m \times m$ orthonormal with first row equal to $\boldsymbol{\alpha}$. so that $\mathbf{H}\boldsymbol{\alpha} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$. Since the distribution of \mathbf{Y} and $\mathbf{Y}\mathbf{H}$ are identical we can without loss of generality set $\boldsymbol{\alpha} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$. Partition $\mathbf{Y} = [\mathbf{Y}_1, \mathbf{Y}_2]$ where \mathbf{Y}_1 is $n \times 1$, \mathbf{Y}_2 is $n \times (m-1)$, and they are independent. Then

$$\begin{aligned} (\boldsymbol{\alpha}'\mathbf{W}^{-1}\boldsymbol{\alpha})^{-1} &= \left(\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} \mathbf{Y}_1'\mathbf{Y}_1 & \mathbf{Y}_1'\mathbf{Y}_2 \\ \mathbf{Y}_2'\mathbf{Y}_1 & \mathbf{Y}_2'\mathbf{Y}_2 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right)^{-1} \\ &= \mathbf{Y}_1'\mathbf{Y}_1 - \mathbf{Y}_1'\mathbf{Y}_2(\mathbf{Y}_2'\mathbf{Y}_2)^{-1}\mathbf{Y}_2'\mathbf{Y}_1 \\ &= \mathbf{Y}_1'\mathbf{M}_2\mathbf{Y}_1 \sim \chi_{n-(m-1)}^2 \end{aligned}$$

where $\mathbf{M}_2 = \mathbf{I}_{m-1} - \mathbf{Y}_2(\mathbf{Y}_2'\mathbf{Y}_2)^{-1}\mathbf{Y}_2'$. The final distributional equality holds conditional on \mathbf{Y}_2 by the same argument in the proof of Theorem 5.7. Since this does not depend on \mathbf{Y}_2 it is the unconditional distribution as well. This establishes the stated result.

To test hypotheses about $\boldsymbol{\mu}$ a classical statistic is known as **Hotelling's** T^2 :

$$T^2 = n(\bar{\mathbf{Y}} - \boldsymbol{\mu})'\hat{\boldsymbol{\Sigma}}^{-1}(\bar{\mathbf{Y}} - \boldsymbol{\mu}).$$

Theorem 11.12 If $Y \sim N(\mu, \Sigma)$ then

$$T^2 \sim \frac{m}{(n-m)(n-1)} F(m, n-m)$$

a scaled F distribution.

To prove this recall that \bar{Y} is independent of $\hat{\Sigma}$. Apply Theorem 11.11 with $\alpha = \bar{Y} - \mu$. Conditional on \bar{Y} and using the fact that $\hat{\Sigma} \sim W_m(n-1, \frac{1}{n-1}\Sigma)$,

$$\begin{aligned} \frac{n}{T^2} &= \left((\bar{Y} - \mu)' \hat{\Sigma}^{-1} (\bar{Y} - \mu) \right)^{-1} \\ &\sim \frac{\chi_{n-1-m+1}^2}{\left(\bar{Y} - \mu \right)' \left(\frac{1}{n-1} \Sigma \right)^{-1} (\bar{Y} - \mu)} \\ &\sim n(n-1) \frac{\chi_{n-m}^2}{\chi_m^2}. \end{aligned}$$

Since the two chi-square variables are independent, this is the stated result.

A very interesting property of this result is that the T^2 statistic is a multivariate quadratic form in normal random variables, yet it has the exact F distribution.

11.18 Exercises

Exercise 11.1 Show (11.10) when the errors are conditionally homoskedastic (11.8).

Exercise 11.2 Show (11.11) when the regressors are common across equations $X_j = X$.

Exercise 11.3 Show (11.12) when the regressors are common across equations $X_j = X$ and the errors are conditionally homoskedastic (11.8).

Exercise 11.4 Prove Theorem 11.1.

Exercise 11.5 Show (11.13) when the regressors are common across equations $X_j = X$.

Exercise 11.6 Show (11.14) when the regressors are common across equations $X_j = X$ and the errors are conditionally homoskedastic (11.8).

Exercise 11.7 Prove Theorem 11.2.

Exercise 11.8 Prove Theorem 11.3.

Exercise 11.9 Show that (11.16) follows from the steps described.

Exercise 11.10 Show that (11.17) follows from the steps described.

Exercise 11.11 Prove Theorem 11.4.

Exercise 11.12 Prove Theorem 11.5.

Hint: First, show that it is sufficient to show that

$$\mathbb{E} \left[\overline{X}' \overline{X} \right] \left(\mathbb{E} \left[\overline{X}' \Sigma^{-1} \overline{X} \right] \right)^{-1} \mathbb{E} \left[\overline{X}' \overline{X} \right] \leq \mathbb{E} \left[\overline{X}' \Sigma \overline{X} \right].$$

Second, rewrite this equation using the transformations $U = \Sigma^{1/2} \overline{X}$ and $V = \Sigma^{1/2} \overline{X}$, and then apply the matrix Cauchy-Schwarz inequality (B.33).

Exercise 11.13 Prove Theorem 11.6.

Exercise 11.14 Take the model

$$\begin{aligned} Y &= \pi' \beta + e \\ \pi &= \mathbb{E}[X | Z] = \Gamma' Z \\ \mathbb{E}[e | Z] &= 0 \end{aligned}$$

where Y is scalar, X is a k vector and Z is an ℓ vector. β and π are $k \times 1$ and Γ is $\ell \times k$. The sample is $(Y_i, X_i, Z_i : i = 1, \dots, n)$ with π_i unobserved.

Consider the estimator $\hat{\beta}$ for β by OLS of Y on $\hat{\pi} = \hat{\Gamma}' Z$ where $\hat{\Gamma}$ is the OLS coefficient from the multivariate regression of X on Z .

- Show that $\hat{\beta}$ is consistent for β .
- Find the asymptotic distribution $\sqrt{n}(\hat{\beta} - \beta)$ as $n \rightarrow \infty$ assuming that $\beta = 0$.
- Why is the assumption $\beta = 0$ an important simplifying condition in part (b)?
- Using the result in (c) construct an appropriate asymptotic test for the hypothesis $\mathbb{H}_0 : \beta = 0$.

Exercise 11.15 The observations are i.i.d., $(Y_{1i}, Y_{2i}, X_i : i = 1, \dots, n)$. The dependent variables Y_1 and Y_2 are real-valued. The regressor X is a k -vector. The model is the two-equation system

$$\begin{aligned} Y_1 &= X' \beta_1 + e_1 \\ \mathbb{E}[X e_1] &= 0 \\ Y_2 &= X' \beta_2 + e_2 \\ \mathbb{E}[X e_2] &= 0. \end{aligned}$$

- What are the appropriate estimators $\hat{\beta}_1$ and $\hat{\beta}_2$ for β_1 and β_2 ?
- Find the joint asymptotic distribution of $\hat{\beta}_1$ and $\hat{\beta}_2$.
- Describe a test for $\mathbb{H}_0 : \beta_1 = \beta_2$.