

## 09 Simultaneous Equations

Andrius Buteikis, andrius.buteikis@mif.vu.lt  
<http://web.vu.lt/mif/a.buteikis/>

# Introduction

Most economic models involve more than one equation. We know how to estimate each equation in, say, supply and demand models, but given the links between the equations, we might reasonably ask whether it is possible to estimate the equations *jointly* (it turns out that it is possible). Sometimes this improves the efficiency of the estimators, sometimes it has no effect on their efficiency, and sometimes it makes things worse.

There are two useful types of relationships between equations.

First, there may be a relationship between the error terms of the model.

Second, there may be relationships between the parameters of the model themselves, either because two parameters take the same value or because some more general mathematical relationship exists between two or more parameters.

## Seemingly Unrelated Regression (SUR)

In the example below, we consider five firms: GM, Chrysler, GE, Westinghouse, and U.S. Steel. The data consist of time series of 20 yearly observations for these firms and three variables:

- ▶  $I_{it}$  - gross investment;
- ▶  $F_{it}$  - market value of the firm at the end of the previous year;
- ▶  $C_{it}$  - value of the stock of plant and equipment at the end of the previous year.

where  $i = 1, \dots, 5$ ,  $t = 1935, \dots, 1954$ . At any moment  $t$  we can write five equations:

$$\left\{ \begin{array}{l} I_{1t} = \beta_{10} + \beta_{11}F_{1t} + \beta_{12}C_{1t} + \epsilon_{1t} \\ \dots\dots\dots\dots\dots\dots\dots \\ I_{5t} = \beta_{50} + \beta_{51}F_{5t} + \beta_{52}C_{5t} + \epsilon_{5t} \end{array} \right.$$

Each equation can be estimated individually taking  $t = 1935, \dots, 1954$  (each firm performs its own investment policy, therefore we can treat these equations as unrelated).

On the other hand, all the economic activities take place in the same economic environment, therefore the five shocks  $\epsilon_{1t}, \dots, \epsilon_{5t}$  can be correlated. The procedure which takes into account this correlation is called **SUR**. Generally, it differs from OLS except for the two cases:

1. The equations are really uncorrelated, i.e.,  $cov(\epsilon_{it}, \epsilon_{js}) = 0, t \neq s$  (and also  $cov(\epsilon_{it}, \epsilon_{js}) = (\delta_{ij})$  does not depend on  $t$ ).
2. All the equations in the above equation system have the same explanatory variables on the right-hand-side, i.e.,

$$F_{1t} = \dots = F_{5t} = F_t \quad \text{etc (this is not true in our case).}$$

The difference between two estimating methods is hardly noticeable, probably, because of a small correlation between the errors. Generally, using SUR to jointly estimate the equations of the system, allowing for correlation between the errors of the equations, will **improve the efficiency** of the estimation, but usually not much.

## Potential Drawbacks of estimating SUR instead of individual OLS

There is one potential problem with simultaneous equations, which is that it requires the Gauss-Markov assumptions to be true for *all* equations. Suppose that the Gauss-Markov assumptions are true for one equation but not for another.

For example, one equation might have an omitted right-hand-side variable or an endogenous one. Then estimating by SUR will generally no longer be unbiased or consistent for *any* of the equations.

In such a case, OLS would remain unbiased and consistent for those equations for which the Gauss Markov assumptions help.

Estimating equation by equation has the advantage that, if there is a problem with one equation, the problem is limited to that equation and cannot spill over to the estimates of the parameters of the other equations.

Up to now, we have not made use of the economic connection between the equations in the system. We have allowed for the error terms to be correlated, and we have some economic ideas about why the errors would be correlated, but the errors might well be correlated by coincidence even if there was no economic link between the equations at all.

We may do better if we can use economic theory to suggest direct links between the parameters of the equations of the system. for example, conditions

$$\beta_2 = \gamma_2 (= \delta)$$

is called a *cross-equation restriction*. If we use OLS to estimate individually both equations from the above **system**, we have to minimize

$$RSS = RSS_1(\beta_0, \beta_1, \delta, \beta_3) + RSS_2(\gamma_0, \gamma_1, \delta, \gamma_3)$$

In this case where we use SUR to estimate the parameters, we have to generalize RSS and to include the effects of the correlation of the error terms.

# Multiple Equations with Endogenous Right-hand-side Variables

SUR is a useful technique for models that can be estimated by least squares.

However, it cannot be used if the Gauss-Markov assumptions are not satisfied.

In particular, if the equations contain **endogenous right-hand-side variables**, SUR will be biased and inconsistent.

Any time we have two equations solving for the values of two variables, such as the supply and demand model, there will be endogenous right-hand-side variables and SUR will not be appropriate.

# Simultaneous Estimation via Three-Stage Least Squares

Fortunately, we can simultaneously estimate equations by **two-stage least squares** in exactly the same way that we can simultaneously estimate them by ordinary least squares. Doing so requires a **three-step process**:

1. Regress each endogenous variable on all exogenous variables in the system of equations, and calculate predicted values for the endogenous variables.
2. Estimate the structural equations by least squares, replacing the endogenous right-hand-side variables with their predicted values from Step 1.
3. Calculate the estimated variances and covariances of the residuals from Step 2, and re-estimate the structural equations using the SUR method.

This technique, known as *three-stage least squares*, is the instrumental variables equivalent to SUR.



It has the same general relationship to two-stage least squares that SUR has to OLS. Its advantage is that it will be more efficient than two-stage least squares for large samples, as long as the right-hand-side variables of the equations are not the same in all equations.

It is not unbiased - but two-stage least squares is not unbiased either - so that is not a disadvantage of three-stage least squares.

Its main disadvantage is that, as with SUR, simultaneous estimation permits a violation of the Gauss-Markov assumptions in one of the equations to spread to the other equations

We have already discussed the problem of the endogenous right-hand-side variables and we presented three methods to deal with it.

The first one was a rather cumbersome indirect least squares method, the other were two - and three-stages least square methods.

Now we shall briefly describe two more methods.

The single equation or limited-information methods (specifically, the *limited information maximum likelihood* method) was introduced in 1949 and was popular until the advent of 2SLS. Computationally it is rather complicated but if the equation under consideration is exactly identified, then LIML and 2SLS give identical estimates.

To estimate the coefficients of an equation, LML uses the information of that equation only. In contrast, in system or *full-information* methods we use information on the restrictions on all equations.

## Example

In 1950, L. Klein proposed the dynamic model of macroeconomics which was later called Klein Model 1. It is described by the following system:

$$\begin{cases} C_t = \alpha_0 + \alpha_1 P_t + \alpha_2 P_{t-1} + \alpha_3 (W_t^P + W_t^G) + \epsilon_{1t} & \text{(consumption)} \\ I_t = \beta_0 + \beta_1 P_t + \beta_2 P_{t-1} + \beta_3 K_{t-1} + \epsilon_{2t} & \text{(investment)} \\ W_t^P = \gamma_0 + \gamma_1 X_t + \gamma_2 X_{t-1} + \gamma_3 A_t + \epsilon_{3t} & \text{(demand for labor)} \\ X_t = C_t + I_t + G_t & \text{(equilibrium demand)} \\ P_t = X_t - T_t - W_t^P & \text{(private sector revenue)} \\ K_t = K_{t-1} + I_t & \text{(capital)} \end{cases}$$

where

- ▶  $C$  - consumption expenditure;  $I$  - investment expenditure;
- ▶  $G$  - government expenditure;  $P$  - profits;
- ▶  $W^P$  - private wage bill;  $W^G$  - government wage bill;
- ▶  $K$  - capital stock;  $A$  - time trend.
- ▶  $T$  - taxes;  $X$  - income after tax;

In the preceding model, the left-hand-side variables  $C$ ,  $I$ ,  $W$ ,  $Y$ ,  $P$ , and  $K$  are treated as jointly dependent, or endogenous, variables,

$G$ ,  $T$ ,  $W^g$ , and  $A$  as exogenous, and the variables  $P_{t-1}$ ,  $K_{t-1}$ , and  $X_{t-1}$  are treated as predetermined.

In all, there are six equations (including the three identities) to study the interdependence of six endogenous variables.

Note that because of the interdependence among the endogenous variables, in general they are not independent of the stochastic disturbance terms, which therefore makes it inappropriate to apply the method of OLS to an individual equation in the system (estimators thus obtained are inconsistent, they do not converge to their true population values even when the sample size is very large).

```

suppressPackageStartupMessages({
  library(readxl)
  require(systemfit)
  require(AER)
})
data("KleinI", package = "systemfit")

```

Specifying the equations in R (while coefficient restrictions can be specified, value restrictions/identities - cannot, so we do not specify them here):

```

eqConsump <- consump ~ corpProf + corpProfLag + wages
eqInvest   <- invest  ~ corpProf + corpProfLag + capitalLag
eqPrivWage <- privWage ~ gnp + gnpLag + trend
inst       <- ~ govExp + taxes + govWage + trend +
              capitalLag + corpProfLag + gnpLag

eq.sys <- list( Consumption = eqConsump,
                Investment   = eqInvest,
                PrivateWages = eqPrivWage )

```

We can estimate the equation system via a variety of estimation techniques:

- ▶ OLS (ignoring the equation system structure - equivalent to single-equation estimation)

```
k.OLS <- systemfit(eq.sys, data = KleinI, method = "OLS")
```

- ▶ 2SLS (equivalent to single-equation instrumental variable estimation)

```
k.2SLS <- systemfit(eq.sys, data = KleinI,  
                    inst = inst, method = "2SLS")
```

- ▶ SUR (residuals are correlated throughout the equations, not accounted for predictor endogeneity problems)

```
k.SUR <- systemfit(eq.sys, data = KleinI, method = "SUR")
```

- ▶ 3SLS (takes into account both the contemporaneous residual correlation and the predictor endogeneity)

```
k.3SLS <- systemfit(eq.sys, data = KleinI,  
                    inst = inst, method = "3SLS")
```

## Estimation method coefficient comparison:

```
cof <- cbind(coef(k.OLS), coef(k.2SLS),  
             coef(k.SUR), coef(k.3SLS))  
colnames(cof) <- c("OLS", "2SLS", "SUR", "3SLS")  
round(cof, 4)
```

##	OLS	2SLS	SUR	3SLS
## Consumption_(Intercept)	16.2366	16.5548	15.9805	16.4408
## Consumption_corpProf	0.1929	0.0173	0.2302	0.1249
## Consumption_corpProfLag	0.0899	0.2162	0.0673	0.1631
## Consumption_wages	0.7962	0.8102	0.7962	0.7901
## Investment_(Intercept)	10.1258	20.2782	12.9293	28.1778
## Investment_corpProf	0.4796	0.1502	0.4429	-0.0131
## Investment_corpProfLag	0.3330	0.6159	0.3655	0.7557
## Investment_capitalLag	-0.1118	-0.1578	-0.1253	-0.1948
## PrivateWages_(Intercept)	1.4970	1.5003	1.6347	1.7972
## PrivateWages_gnp	0.4395	0.4389	0.4098	0.4005
## PrivateWages_gnpLag	0.1461	0.1467	0.1744	0.1813
## PrivateWages_trend	0.1302	0.1304	0.1558	0.1497

Estimation method coefficient p-value comparison:

```
p.val <- cbind(summary(k.OLS)$coefficient[, 4], summary(k.2SLS)$  
              summary(k.SUR)$coefficient[, 4], summary(k.3SLS)$co  
colnames(p.val) <- c("OLS", "2SLS", "SUR", "3SLS")  
round(p.val, 4)
```

##	OLS	2SLS	SUR	3SLS
## Consumption_(Intercept)	0.0000	0.0000	0.0000	0.0000
## Consumption_corpProf	0.0495	0.8966	0.0152	0.3133
## Consumption_corpProfLag	0.3353	0.0874	0.4422	0.1621
## Consumption_wages	0.0000	0.0000	0.0000	0.0000
## Investment_(Intercept)	0.0814	0.0271	0.0269	0.0017
## Investment_corpProf	0.0001	0.4460	0.0002	0.9429
## Investment_corpProfLag	0.0042	0.0034	0.0019	0.0004
## Investment_capitalLag	0.0006	0.0011	0.0002	0.0000
## PrivateWages_(Intercept)	0.2547	0.2558	0.2055	0.1655
## PrivateWages_gnp	0.0000	0.0000	0.0000	0.0000
## PrivateWages_gnpLag	0.0011	0.0034	0.0001	0.0002
## PrivateWages_trend	0.0008	0.0009	0.0001	0.0002



Estimation method coefficient standard error comparison:

```
s.e <- cbind(summary(k.OLS)$coefficient[, 2], summary(k.2SLS)$coefficient[, 2],  
             summary(k.SUR)$coefficient[, 2], summary(k.3SLS)$coefficient[, 2])  
colnames(s.e) <- c("OLS", "2SLS", "SUR", "3SLS")  
round(s.e, 4)
```

##	OLS	2SLS	SUR	3SLS
## Consumption_(Intercept)	1.3027	1.4680	1.2989	1.4499
## Consumption_corpProf	0.0912	0.1312	0.0852	0.1202
## Consumption_corpProfLag	0.0906	0.1192	0.0855	0.1116
## Consumption_wages	0.0399	0.0447	0.0392	0.0422
## Investment_(Intercept)	5.4655	8.3832	5.3364	7.5509
## Investment_corpProf	0.0971	0.1925	0.0957	0.1799
## Investment_corpProfLag	0.1009	0.1809	0.0994	0.1700
## Investment_capitalLag	0.0267	0.0402	0.0261	0.0362
## PrivateWages_(Intercept)	1.2700	1.2757	1.2418	1.2402
## PrivateWages_gnp	0.0324	0.0396	0.0303	0.0354
## PrivateWages_gnpLag	0.0374	0.0432	0.0347	0.0380
## PrivateWages_trend	0.0319	0.0324	0.0307	0.0310

## Summary

The set of economic variables  $Y_1, \dots, Y_k$  is determined through a market equilibrium mechanism and we want to analyze the structure of relationships that determines the equilibrium. Suppose that

$$\vec{Y} = (Y_1, \dots, Y_n)$$

is a vector consisting of  $n$  economic variables, among which there exist  $n$  relationships that determine the equilibrium levels of the variables. We also suppose that there exist  $x$  variables

$$\vec{Z} = (Z_1, \dots, Z_x)$$

that are independent of the economic relations but affect the equilibrium.

The variables  $\hat{Y}$  are called **endogenous** variables, and the  $\hat{Z}$  are called **exogenous** variables.

If we assume linear relationships among them, we have an expression such as

$$\vec{Y} = B\vec{Y} + \Gamma\vec{Z} + \vec{u}$$

where  $B$  and  $\Gamma$  are matrices with constant coefficients and  $\vec{u}$  is a vector of disturbances or errors.

The mentioned system is called the linear *structural equation system* and is a system of simultaneous equations. By solving the equations formally, we get the so-called *reduced* form

$$\vec{Y} = \Pi\vec{Z} + \vec{v}$$

where

$$\Pi = (I - B)^{-1}\Gamma,$$

$$\vec{v} = (I - B)^{-1}\vec{u}$$

The relation of  $\vec{Y}$  to  $\vec{Z}$  is determined through the reduced form, and if we have enough data on  $\vec{Y}$  and  $\vec{Z}$ , we can estimate  $\Pi$ .

The problem of identification is to decide whether we can determine the unknown parameters in  $B$  and  $\Gamma$  uniquely from the parameters in the reduced form.

A necessary condition for the parameters in **one** of the equations in the first system ( $\vec{Y} = B\vec{Y} + \Gamma\vec{Z} + \vec{u}$ ) to be identifiable is that the number of unknown variables in the equation not be greater than  $x + 1$ .

If it is exactly equal to  $x + 1$ , the equation is said to be *just identified*, and if it is less than  $x + 1$ , the equation is said to be *overidentified*.

- ▶ If all the equations in the system are **just identified**, for arbitrary  $\Pi$  there exist unique  $B$  and  $\Gamma$  that satisfy

$$\Pi = (I - B)^{-1}\Gamma$$

Therefore, if we denote the least squares estimator of  $\Pi$  by  $\hat{\Pi}$ , we can estimate  $B$  and  $\Gamma$  from the equation

$$(I - \hat{B})\hat{\Pi} = \hat{\Gamma}$$

This procedure is called the *indirect least squares* method and is equivalent to the maximum likelihood method if we assume normality for  $\vec{u}$

- ▶ When some of the equations are **overidentified**, the estimation problem becomes complicated. Three kinds of procedures have been proposed:
  1. Full system methods
  2. Single equation methods
  3. Subsystem methods

In full system methods all the parameters are considered simultaneously, and if normality is assumed, the maximum likelihood estimator can be obtained by minimizing

$$|(\vec{Y} - \Pi\vec{Z})(\vec{Y} - \Pi\vec{Z})'|$$

Since it is usually difficult to compute the maximum likelihood estimator, a simpler, but asymptotically equivalent, *three stage least squares* method has been proposed.

The **single equation** methods and the **subsystem** methods take into consideration only the information about the parameters in one equation or in a subset of the equations, and estimate the parameters in each equation separately.

There is a **single equation** method, called the *limited information maximum likelihood* method, based on the maximum likelihood approach, and also a *two-stage squares* method, which estimated  $\Pi$  first by least squares, computes  $\hat{\vec{Y}} = \hat{\Pi}\vec{Z}$ , and then applies the least squares method to the model  $\vec{Y} = B\hat{\vec{Y}} + \Gamma\vec{Z} + \vec{u}$ . These two and also some others are asymptotically equivalent.

Among asymptotically equivalent classes of estimators corresponding to different information structures it has been established that the maximum likelihood estimators have asymptotically higher-order efficiency than other estimators, and Monte Carlo and numerical studies show that they are in most cases better than others if properly adjusted for the biases.

In many simultaneous equation models which have been applied to actual macroeconomic data, the values of endogenous variables obtained in the past appear on the right-hand sides of equations  $\vec{Y} = B\vec{Y} + \Gamma\vec{Z} + \vec{u}$ .

Such variables are called *lagged variables*, and they can be treated, at least in the asymptotic theory of inference, as though they were exogenous.

Hence exogenous variables and lagged endogenous variables are jointly called *predetermined variables*. When many lagged variables appear over many time periods and when some structure among the coefficients of those lagged variables can be assumed, such a model is called a *distributed lag model*.

# Summary

## Simultaneous Equations Model

The model:

- ▶ Consists of a set of equations.
- ▶ Has two, or more, dependent variables;

Simultaneous equations models require special statistical methods for parameter estimation as the least squares estimation is not appropriate for these models.



For example, price and quantity are determined by the *interaction* of the two equations - one for supply, and one for demand:

$$\text{Demand: } Q = \alpha_1 P + \alpha_2 X + \epsilon_{demand}$$

$$\text{Supply: } Q = \beta_1 P + \epsilon_{supply}$$

where:

- ▶  $Q$  - the quantity demanded;
- ▶  $P$  - the price;
- ▶  $X$  - the income;
- ▶  $Cov(\epsilon_{demand}, \epsilon_{supply}) = 0$ .

**Two** equations are used to describe the supply and demand equilibrium; **Two** equilibrium values (price  $P^*$  and quantity  $Q^*$ ) are determined at the **same time**.

In this case,  $P$  and  $Q$  are the **endogenous variables** (and are therefore random variables) since their values are determined within the specified system. The income variable  $X$  is determined **outside the system** and, hence, is called an **exogenous** variable.

Note that we usually assume that the left-hand-side **explanatory** variables are fixed (i.e. non-random). Because  $P$  is on the left and is an endogenous variable and is therefore a random variable, this means that  $P$  and  $Q$  are correlated which means that  $\epsilon_{supply}$  is correlated with  $P$ . To show this, we need to rewrite the equation system.

## Reduced-Form Equations

The previously mentioned equations can be solved to express the endogenous variables  $P$  and  $Q$  as functions of the remaining exogenous variables (in our case it is only one -  $X$ ). This reformulated model is called the **reduced form** of the structural equation system.

- ▶ To solve for  $P$ , we set  $Q$  in the demand and supply equations to be equal:

$$\beta_1 P + \epsilon_{supply} = \alpha_1 P + \alpha_2 X + \epsilon_{demand}$$

Then we solve for  $P$ :

$$\begin{aligned} P &= \frac{\alpha_2}{\beta_1 - \alpha_1} X + \frac{\epsilon_{demand} - \epsilon_{supply}}{\beta_1 - \alpha_1} \\ &= \pi_1 X + v_1 \end{aligned}$$

- ▶ To solve for  $Q$ , we substitute the new expression of  $P$  into either demand or supply equations. Let's substitute it into the supply eq. to get:

$$\begin{aligned} Q &= \beta_1 P + \epsilon_{supply} \\ &= \beta_1 \left[ \frac{\alpha_2}{\beta_1 - \alpha_1} X + \frac{\epsilon_{demand} - \epsilon_{supply}}{\beta_1 - \alpha_1} \right] + \epsilon_{supply} \\ &= \frac{\beta_1 \alpha_2}{\beta_1 - \alpha_1} X + \frac{\beta_1 \epsilon_{demand} - \alpha_1 \epsilon_{supply}}{\beta_1 - \alpha_1} \\ &= \pi_2 X + v_2 \end{aligned}$$

The parameters  $\pi_1$  and  $\pi_2$  are called the **reduced-form parameters**.  
The error terms  $v_1$  and  $v_2$  are called the **reduced-form errors**.

The reduced form equations can be estimated consistently by OLS to get  $\hat{\pi}_1$  and  $\hat{\pi}_2$ . Because  $X$  is determined outside of the equation, it is not correlated with  $v_1$  and  $v_2$ .

What if we want to estimate the original  $\alpha_1$  and  $\beta_1$  with the original system unchanged?

- ▶ OLS produces incorrect estimates because as mentioned:

$$\begin{aligned} \text{Cov}(P, \epsilon_{supply}) &= \text{Cov}(\pi_1 X + v_1, \epsilon_{supply}) \\ &= \text{Cov}\left(\frac{\epsilon_{demand} - \epsilon_{supply}}{\beta_1 - \alpha_1}, \epsilon_{supply}\right) \\ &= \frac{1}{\beta_1 - \alpha_1} \text{Cov}(-\epsilon_{supply}, \epsilon_{supply}) \\ &= \frac{-\sigma_{\epsilon, supply}^2}{\beta_1 - \alpha_1} < 0 \end{aligned}$$

Note that for  $\beta_1 P + \epsilon_{supply} = \alpha_1 P + \alpha_2 X + \epsilon_{demand}$  to hold, we need  $\beta_1 \geq \alpha_1$ , because  $X \geq 0$ .

The reduced-form equations are important for economic analysis. These equations relate the equilibrium values of the endogenous variables to the exogenous variables. This lets us analyse what magnitude an increase in  $X$  could have to price, and what would be the new equilibrium of the adjusted market in terms of  $P$  and  $Q$ .

The estimated reduced-form equations can be used to predict values of equilibrium price and quantity for different levels of income.

## Identification

In a system of  $M$  simultaneous equations, which jointly determine  $M$  endogeneous variables, at least  $M - 1$  variables must be excluded from one of the of the equations for estimation of its parameters to be possible. Then, such an equation is said to be **identified** and its parameters can be estimated consistently.

If fewer than  $M - 1$  variables are omitted from an equation, then it is said to be unidentified and its parameters cannot be consistently estimated.

In our example there are  $M = 2$  equations, so we need  $M - 1 = 1$  variable to be omitted from an equation to identify it. There are a total of three variables:  $P$ ,  $Q$  and  $X$ .

The demand equation does not have any variables omitted - it is **unidentified** and its parameters cannot be estimated consistently.

The supply equation does not have the variable  $X$ , so it is **identified** and its parameters can be estimated.

The number of instrumental variables required for estimation of an equation within a simultaneous equations model is equal to the number of right-hand-side endogenous variables.

## Two-Stage Least Squares Estimation

Assume that we have  $M$  endogeneous variables  $Y_1, \dots, Y_M$ ,  $K$  exogeneous variables  $X_1, \dots, X_K$  and suppose that **the first** structural equation is:

$$Y_1 = \alpha_2 Y_2 + \alpha_3 Y_3 + \beta_1 X_1 + \beta_2 X_2 + e_1$$

If the equation is **identified**, then its parameters can be estimated in **two steps**:

1. Estimate the parameters of the reduced form equations by OLS:

$$Y_2 = \pi_{12}X_1 + \pi_{22}X_2 + \dots + \pi_{K2}X_K + v_2$$

$$Y_3 = \pi_{13}X_1 + \pi_{23}X_2 + \dots + \pi_{K3}X_K + v_3$$

and obtain the predicted values:

$$\hat{Y}_2 = \hat{\pi}_{12}X_1 + \hat{\pi}_{22}X_2 + \dots + \hat{\pi}_{K2}X_K$$

$$\hat{Y}_3 = \hat{\pi}_{13}X_1 + \hat{\pi}_{23}X_2 + \dots + \hat{\pi}_{K3}X_K$$

2. Replace the endogeneous variables  $Y_2$  and  $Y_3$  on the right-hand-side of the structural equation of  $Y_1$  with their predicted values:

$$Y_1 = \alpha_2 \hat{Y}_2 + \alpha_3 \hat{Y}_3 + \beta_1 X_1 + \beta_2 X_2 + e_1^*$$

Estimate the parameters of this equation by OLS.

```
install.packages("devtools")
library(devtools)
install_git("https://github.com/ccolonescu/PoEdata")
```

## Example

We will try to predict the supply and demand for truffles.

```
data("truffles", package="PoEdata")
head(truffles)
```

##	p	q	ps	di	pf
## 1	29.64	19.89	19.97	2.103	10.52
## 2	40.23	13.04	18.04	2.043	19.67
## 3	34.71	19.61	22.36	1.870	13.74
## 4	41.43	17.13	20.87	1.525	17.95
## 5	53.37	22.55	19.79	2.709	13.71
## 6	38.52	6.37	15.98	2.489	24.95



Here  $p$  - the price of truffles,  $q$  - is the quantity of truffles traded in a particular French market,  $ps$  - price of substitute for real truffles,  $di$  - per capita monthly disposable income,  $pf$  - price factor of production (in this case - hourly rental price of truffle-pigs used in the search process).

Economic theory says that the price of a factor of production should affect supply but not demand, and that the price of substitute goods and income should affect demand and not supply. The specifications we use are based on the microeconomic theory of supply and demand. So, our supply and demand model for truffles is:

$$\text{Supply: } Q_i = \beta_1 + \beta_2 P_i + \beta_3 PF_i + \epsilon_{si}$$

$$\text{Demand: } Q_i = \alpha_1 + \alpha_2 P_i + \alpha_3 PS_i + \alpha_4 DI_i + \epsilon_{di}$$

Here  $i$  indicates a French market  $i$ . In this model we assume that  $P$  and  $Q$  are endogenous variables and  $PS$ ,  $DI$ ,  $PF$  and the *intercept* are exogenous variables.

In this case we have  $M = 2$  equations. So, at least  $M - 1 = 1$  variables must be excluded from each equation in order for it to be identified. In the demand equation -  $PF$  is not included; in the supply equation -  $PS$  and  $DI$  are not included. So, the system is **identified**.

The reduced-form equations express the endogenous variables  $P$  and  $Q$  as a function of the exogenous variables  $PS$ ,  $DI$ ,  $OF$  and the intercept:

$$Q_i = \pi_{11} + \pi_{21}PS_i + \pi_{31}DI_i + \pi_{41}PF_i + v_{i1}$$

$$P_i = \pi_{12} + \pi_{22}PS_i + \pi_{32}DI_i + \pi_{42}PF_i + v_{i2}$$

We can estimate these equations by least squares since the right-hand-side variables are exogenous and uncorrelated with the random errors  $v_{i1}$  and  $v_{i2}$ .

```
#Estimate reduced-form parameters
```

```
q.lm <- lm(q ~ ps + di + pf, data = truffles)  
p.lm <- lm(p ~ ps + di + pf, data = truffles)  
round(summary(q.lm)$coefficients, 4)
```

```
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept)   7.8951     3.2434   2.4342  0.0221  
## ps            0.6564     0.1425   4.6051  0.0001  
## di            2.1672     0.7005   3.0938  0.0047  
## pf           -0.5070     0.1213  -4.1809  0.0003
```

```
round(summary(p.lm)$coefficients, 4)
```

```
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -32.5124     7.9842  -4.0721  4e-04  
## ps            1.7081     0.3509   4.8682  0e+00  
## di            7.6025     1.7243   4.4089  2e-04  
## pf            1.3539     0.2985   4.5356  1e-04
```

The reduced-form equations are used to obtain  $\hat{P}$ :

$$\hat{P} = -32.51 + 1.71PS + 7.6DI + 1.36PF$$

that will be used in place of  $P$  on the right-hand side of the supply and demand equations in the second stage of two-stage least squares.

```
truffles$p_hat <- p.lm$fitted.values  
demand.lm <- lm(q ~ p_hat + ps + di, data = truffles)  
supply.lm <- lm(q ~ p_hat + pf, data = truffles)
```

```
round(summary(demand.lm)$coefficients, 4)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	-4.2795	3.0138	-1.4199	0.1675
##	p_hat	-0.3745	0.0896	-4.1809	0.0003
##	ps	1.2960	0.1931	6.7119	0.0000
##	di	5.0140	1.2414	4.0389	0.0004

```
round(summary(supply.lm )$coefficients, 4)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	20.0328	2.1657	9.2500	0
##	p_hat	0.3380	0.0441	7.6599	0
##	pf	-1.0009	0.1461	-6.8496	0

Note that by applying two least squares regressions using ordinary least squares regression software, the standard errors and t-values reported in the second regression **are not correct for the 2SLS estimator**.

To get around this, we should use specialized functions for 2SLS estimation, like `tsls()` from `sem` package, or the `ivreg()` function from `AER` package.

```
suppressPackageStartupMessages({library(sem)})
demand <- tsls(q ~ p + ps + di, ~ ps + di + pf, data = truffles)
supply <- tsls(q ~ p + pf, ~ ps + di + pf, data = truffles)
round(summary(demand)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-4.2795	5.5439	-0.7719	0.4471
## p	-0.3745	0.1648	-2.2729	0.0315
## ps	1.2960	0.3552	3.6488	0.0012
## di	5.0140	2.2836	2.1957	0.0372

```
round(summary(supply)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	20.0328	1.2231	16.3785	0
## p	0.3380	0.0249	13.5629	0
## pf	-1.0009	0.0825	-12.1281	0

```
suppressPackageStartupMessages({library(AER)})
demand <- ivreg(q ~ p + ps + di | ps + di + pf, data = truffles)
supply <- ivreg(q ~ p + pf | ps + di + pf, data = truffles)
round(summary(demand)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-4.2795	5.5439	-0.7719	0.4471
## p	-0.3745	0.1648	-2.2729	0.0315
## ps	1.2960	0.3552	3.6488	0.0012
## di	5.0140	2.2836	2.1957	0.0372

```
round(summary(supply)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	20.0328	1.2231	16.3785	0
## p	0.3380	0.0249	13.5629	0
## pf	-1.0009	0.0825	-12.1281	0

The estimated equations are:

$$\widehat{\text{Supply}} = 20.03 + 0.34\hat{P} - 1.00PF$$
$$\widehat{\text{Demand}} = -4.28 - 0.37\hat{P} + 1.30PS + 5.01DI$$

Note that both `tsts()` and `ivreg()` produce the same results, but the `t` value and `p` value of the estimated coefficients in `lm()` are different!

The most common technique of solving for simultaneous equation models is a technique called two-staged least squares. This method transforms a set of simultaneous equations into functional forms that use the endogenous variables as a function of the system's exogenous variables. You can then use least squares to get the estimators for the reduced-form equations. The final step is to plug one of the fitted values into the right-hand side of one of your structural equations to get the correct estimates of your equations.



## Seemingly Unrelated Regressions

Lets say we have two equations of different firm investments:

$$INV_{1,t} = \beta_{11} + \beta_{12}V_{1,t} + \beta_{13}K_{1,t} + e_{1,t}$$

$$INV_{2,t} = \beta_{21} + \beta_{22}V_{2,t} + \beta_{23}K_{2,t} + e_{2,t}$$

Where  $INV_{i,t}$  is the gross investment for firm  $i$  at time  $t$ ;  $V_{i,t}$  is the stock market value of firm  $i$  at time  $t$  and  $K_{i,t}$  is the actual capital stock of firm  $i$  at time  $t$ .

The two investment equations appeared unrelated and we may estimate them separately.

```
suppressPackageStartupMessages({library(plm)})  
data("grunfeld2", package="PoEdata")  
grun <- pdata.frame(grunfeld2, index = c("firm", "year"))  
head(grun)
```

```
##           inv       v      k firm year  
## 1-1935 33.1 1170.6  97.8    1 1935  
## 1-1936 45.0 2015.8 104.4    1 1936
```

```
inv1_lm <- lm(inv ~ v + k, data = grun[grun$firm == 1, ])  
inv2_lm <- lm(inv ~ v + k, data = grun[grun$firm == 2, ])  
round(summary(inv1_lm)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-9.9563	31.3742	-0.3173	0.7548
## v	0.0266	0.0156	1.7057	0.1063
## k	0.1517	0.0257	5.9015	0.0000

```
round(summary(inv2_lm)$coefficients, 4)
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-0.5094	8.0153	-0.0636	0.9501
## v	0.0529	0.0157	3.3677	0.0037
## k	0.0924	0.0561	1.6472	0.1179

The equations have different coefficients and their error variances:

```
var(resid(inv1_lm))
```

```
## [1] 695.6099
```

```
var(resid(inv2_lm))
```

```
## [1] 93.3281
```

are also different. If additionally, the errors in one equation are **uncorrelated** to the errors in the other equation, then we do indeed have nothing to link the two equations together - combining the data from the two firms brings no gain.

On the other hand, if  $Cov(e_{1,t}, e_{2,t}) \neq 0$ , i.e. if the error terms in the two equations are correlated, then we should employ a different approach.

In the context of this dataset, the errors contain the influence on investment of factors that have been omitted from the equations. Such factors might include:

- ▶ capacity utilization;
- ▶ current and past interest rates;
- ▶ liquidity;
- ▶ the general state of the economy.

If the two firms are similar in many respects (if both of them are operating in the same markets), it is likely that the effects of the omitted factors on investment by one company will be similar to their effect on investment by the other company.

If so, then the error terms will be capturing similar effects and will be correlated.

Adding the **contemporaneous correlation** (correlation between the realizations of two time series variables in the same time period) assumption has the effect of introducing additional information that is not included when we carry out separate least squares estimation of the two equations.

The seemingly unrelated regressions (**SUR**) estimation generalizes the least squares estimation procedure - it estimates the two investment equations jointly, taking into account that:

- ▶ the variances of the error terms are different for the two equations;
- ▶ the errors in the two equations are contemporaneously correlated.

There are **three stages** in the SUR estimation procedure:

1. Estimate the equations separately using OLS.
2. Use the residuals from step 1 to estimate  $Var(e_{1,t})$ ,  $Var(e_{2,t})$ ,  $Cov(e_{1,t}, e_{2,t})$ .
3. Use the estimates from step 2 to estimate the two equations jointly within a generalized least squares (GLS) framework.

In general, if our equation is  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$  and the error variance-covariance matrix is  $\mathbb{E}(\epsilon'\epsilon) = \sigma^2\Omega$ , then:

- ▶ using OLS we estimate the parameters via:

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

- ▶ using GLS we estimate the parameters via:

$$\hat{\beta} = (\mathbf{X}'\Omega^{-1}\mathbf{X})^{-1} \mathbf{X}'\Omega^{-1}\mathbf{Y}$$

```
colnames(grun) <- c("inv", "val", "cap", "firm", "year")
```

We can estimate the parameters via `systemfit()` from the package of the same name:

```
suppressPackageStartupMessages({library("systemfit")})  
grunf.SUR <- systemfit(inv ~ val + cap,  
                       method = "SUR",  
                       data = grun)  
round(summary(grunf.SUR)$coefficients, 4)
```

##		Estimate	Std. Error	t value	Pr(> t )
##	1_(Intercept)	-27.7193	29.3212	-0.9454	0.3577
##	X1_val	0.0383	0.0144	2.6576	0.0166
##	X1_cap	0.1390	0.0250	5.5647	0.0000
##	2_(Intercept)	-1.2520	7.5452	-0.1659	0.8702
##	X2_val	0.0576	0.0145	3.9618	0.0010
##	X2_cap	0.0640	0.0530	1.2062	0.2443

Residual correlation matrix estimate:

```
summary(grunf.SUR)$residCor
```

```
##           1           2
## 1 1.0000000 0.7650429
## 2 0.7650429 1.0000000
```

The SUR estimation procedure is optimal under the contemporaneous correlation assumption, so no standard error adjustment is necessary.

Since the SUR technique utilizes the information on the correlation between the error terms, it is more precise than the OLS - the standard errors of the SUR estimates are lower than those of the OLS.

You should be cautious, when making judgments about precision on the basis of standard errors. Standard errors are themselves estimates; it is possible for a standard error for SUR to be greater than a corresponding least squares standard error even when SUR is a better estimator than least squares.

So, the equations seemed to be unrelated, but the additional information provided by the correlation between the equation errors means that joint GLS estimation is better than single-equation OLS estimation.

There are two situations in which separate least squares estimation is just as good as the SUR technique:

1. when the equation errors are not contemporaneously correlated;
2. when the same explanatory variables (i.e. with the same observation values) appear in each equation.