

Conditional Expectation Function (CEF)

- We begin by thinking about population relationships.
- CEF Decomposition Theorem: Given some outcome Y_i and some covariates X_i there is always a decomposition

$$Y_i = E(Y_i/X_i) + \epsilon_i \quad (1)$$

where

$$E(\epsilon_i/X_i) = 0 \quad (2)$$

- Proof:

$$E(\epsilon_i/X_i) = E[(Y_i - E(Y_i/X_i))/X_i] \quad (3)$$

$$= E(Y_i/X_i) - E[E(Y_i/X_i)/X_i] \quad (4)$$

$$= 0 \quad (5)$$

- The last step uses the *Law of Iterated Expectations*:

$$E(Y) = E[E(Y/X)] \quad (6)$$

where the outer expectation is over X . For example, the average outcome is the weighted average of the average outcome for men and the average outcome for women where the weights are the proportion of each sex in the population.

- The CEF Decomposition Theorem implies that ϵ_i is uncorrelated with any function of X_i .
- Proof: Let $h(X_i)$ be any function of X_i .

$$E[h(X_i)\epsilon_i] = E\{E[(h(X_i)\epsilon_i)/X_i]\} = E\{h(X_i)E(\epsilon_i/X_i)\} = 0 \quad (7)$$

- We refer to the $E(Y_i/X_i)$ as the CEF.

- **Best Predictor:** $E(Y_i/X_i)$ is the Best (Minimum mean squared error – MMSE) predictor of Y_i in that it minimises the function

$$E((Y_i - h(X_i))^2) \quad (8)$$

where $h(X_i)$ is any function of X_i .

- A regression model is a particular choice of function for $E(Y_i/X_i)$.
- Linear regression:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} + \varepsilon \quad (9)$$

- Of course, the linear model may not be correct.

Linear Predictors

- A linear predictor (with only one regressor) takes the form

$$E^*(Y_i/X_i) = \beta_1 + \beta_2 X_{2i} \quad (10)$$

- Suppose we want the Best Linear Predictor (BLP) for Y_i to minimise

$$E((Y_i - E^*(Y_i/X_i))^2) \quad (11)$$

- The solution is

$$\beta_1^* = \mu_Y - \beta_2^* \mu_X \quad (12)$$

$$\beta_2^* = \text{cov}(X_i, Y_i) / \text{Var}(X_i) \quad (13)$$

- In the multivariate case with

$$E^*(Y_i/X_i) = X_i'\beta = \beta_1 + \beta_2 X_{2i} + \dots + \beta_K X_{Ki} \quad (14)$$

, we have

$$\beta^* = E(X_i X_i')^{-1} E(X_i Y_i) \quad (15)$$

- We can see this by taking the first order condition for (11)

$$X_i(Y_i - X_i'\beta) = 0 \quad (16)$$

- The error term $\varepsilon_i = Y_i - E^*(Y_i/X_i)$ satisfies $E(\varepsilon_i X_i) = 0$.
- $E^*(Y/X)$ is the best linear approximation to $E(Y/X)$.
- If $E(Y_i/X_i)$ is linear in X_i , then $E(Y_i/X_i) = E^*(Y_i/X_i)$.

- Will only be the case that $E(\varepsilon_i/X_i) = 0$ if the CEF is linear. However $E(\varepsilon_i X_i) = 0$ in all cases.

Example: Bivariate Normal Distribution

- Assume Z_1 and Z_2 are two standard normal variables and

$$X_1 = \mu_1 + \sigma_1 Z_1 \quad (17)$$

$$X_2 = \mu_2 + \sigma_2(\rho Z_1 + \sqrt{(1 - \rho^2)} Z_2) \quad (18)$$

- Then (X_1, X_2) are bivariate normal with $X_j \sim N(\mu_j, \sigma_j^2)$.
- The covariance between (X_1, X_2) is $\rho\sigma_1\sigma_2$.

- Then using (12) and (13), the BLP

$$E^*(X_2/X_1) = \mu_2 - \beta\mu_1 + \beta X_1 \quad (19)$$

$$= \mu_2 + \beta(X_1 - \mu_1) \quad (20)$$

where

$$\beta = \rho\sigma_2/\sigma_1 \quad (21)$$

- Note that using properties of the Normal distribution,

$$E(X_2/X_1) = \mu_2 + \beta(X_1 - \mu_1) \quad (22)$$

so the CEF is linear in this case.

- This is not true with other distributions.

CEF and Regression Summary

- If the CEF is linear, then the population regression function is exactly it.
- If the CEF is non-linear, the regression function provides the best linear approximator to it.
- The CEF is linear only in special cases such as:

1. Joint Normality of Y and X .
2. Saturated regression models – models with a separate parameter for every possible combination of values that the regressors can take. This occurs when there is a full set of dummy variables and interactions between the dummies.

For example, suppose we have a dummy for female (x_1) and a dummy for white (x_2). The CEF is

$$E(Y/x_1, x_2) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \delta x_1 x_2 \quad (23)$$

We can see that there is a parameter for each possible set of values

$$E(Y/x_1 = 0, x_2 = 0) = \alpha \quad (24)$$

$$E(Y/x_1 = 0, x_2 = 1) = \alpha + \beta_2 \quad (25)$$

$$E(Y/x_1 = 1, x_2 = 0) = \alpha + \beta_1 \quad (26)$$

$$E(Y/x_1 = 1, x_2 = 1) = \alpha + \beta_1 + \beta_2 + \delta \quad (27)$$

Linear Regression

- The regression model is

$$Y_i = X_i' \beta + \varepsilon_i \quad (28)$$

- We assume that $E(\varepsilon_i/X_i) = 0$ if the CEF is linear.
- We assume only that $E(X_i\varepsilon_i) = 0$ if we believe the CEF is nonlinear.
- The population parameters are

$$\beta = E(X_iX_i')^{-1}E(X_iY_i) \quad (29)$$

- Here, Y_i is a scalar and X_i is $K * 1$ where K is the number of X variables.
- In matrix notation,

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, X = \begin{pmatrix} X'_1 \\ X'_2 \\ \vdots \\ X'_n \end{pmatrix}, \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad (30)$$

- So, we can write the model as

$$Y = X\beta + \varepsilon \quad (31)$$

where Y is $n * 1$, X is $n * K$, β is $K * 1$, and ε is $n * 1$.

Best Linear Predictor

- We observe n observations on $\{Y_i, X_i\}$ for $i = 1, \dots, n$.
- We derive the OLS estimator as the BLP as it solves the sample analog of $\min E((Y_i - X_i'\beta)^2)$.
- In fact it minimises $1/n \sum_i (Y_i - X_i'b)^2$.

- In matrix notation, this equals $(Y - Xb)'(Y - Xb) = Y'Y - 2b'X'Y + b'X'Xb$
- FOC: $X'Y - X'Xb = X'(Y - Xb) = 0$.
- So long as $X'X$ is invertible (X is of full rank), $b = (X'X)^{-1}X'Y$

Regression Basics

- The fitted value of a regression

$$\hat{Y} = Xb = X(X'X)^{-1}X'Y = PY \quad (32)$$

where $P = X(X'X)^{-1}X'$.

- The residuals

$$e = Y - Xb \quad (33)$$

$$= Y - X(X'X)^{-1}X'Y \quad (34)$$

$$= (I - P)Y = MY \quad (35)$$

- M and P are symmetric idempotent matrices.
- Any matrix A is idempotent if it is square and $AA = A$.
- P is called the projection matrix that projects Y onto the columns of X to produce the set of fitted values

$$\hat{Y} = Xb = PY \quad (36)$$

- What happens when you project X onto X ?

- M is the residual-generating matrix as $e = MY$. We can easily show that

$$MY = M\varepsilon \quad (37)$$

so although the true errors are unobserved, we can obtain a certain linear combination.

Frisch-Waugh-Lovell Theorem

- Consider the partitioned regression

$$Y = X_1\beta_1 + X_2\beta_2 + \varepsilon \quad (38)$$

- Here X_1 is $n * K_1$ and X_2 is $n * K_2$.

- The FOC are

$$\begin{pmatrix} X_1' \\ X_2' \end{pmatrix} (Y - X_1 b_1 - X_2 b_2) = 0 \quad (39)$$

- Therefore

$$\begin{aligned} X_1'(Y - X_1 b_1 - X_2 b_2) &= 0 \\ X_2'(Y - X_1 b_1 - X_2 b_2) &= 0 \end{aligned} \quad (40)$$

and

$$X_1' X_1 b_1 + X_1' X_2 b_2 = X_1' Y \quad (41)$$

$$X_2' X_1 b_1 + X_2' X_2 b_2 = X_2' Y \quad (42)$$

- Solving these simultaneous equations (you should be able to do this) we get

$$b_1 = (X_1' M_2 X_1)^{-1} X_1' M_2 Y \quad (43)$$

$$b_2 = (X_2' M_1 X_2)^{-1} X_2' M_1 Y \quad (44)$$

where $M_1 = I - X_1(X_1'X_1)^{-1}X_1'$ and $M_2 = I - X_2(X_2'X_2)^{-1}X_2'$.

- Implies can estimate b_1 by regressing M_2Y on M_2X_1 .
- That is, can first regress Y and X_1 on X_2 and then regress residuals on residuals.
- Note that one can also do this by regressing only X_1 on X_2 and then regressing Y on the residuals.

Application to the Intercept

- Consider the case where X_1 is the intercept (a column vector of ones)
 $X_1 = l$.

- Now

$$M_1 = I - l(l'l)^{-1}l' = I - \frac{1}{n}ll' \quad (45)$$

- Note that $\frac{1}{n}ll'$ is an $n * n$ matrix with each element equal to $1/n$.
- Therefore, in this case,

$$M_1X_2 = X_2 - \bar{X}_2 \quad (46)$$

- Implies that regressing Y and X_2 on a constant and then regressing residuals on residuals is the same as taking deviations from means.
- Also implies that if you remove means from all variables, you do not need to include a constant term.

Example: Non-Linear and Linear CEF -- Birth Order and IQ in 5-Child Families (N=10214)

1. X includes a constant and a linear Birth Order variable.

$$\beta_1 = 5.2, \beta_2 = -0.10$$

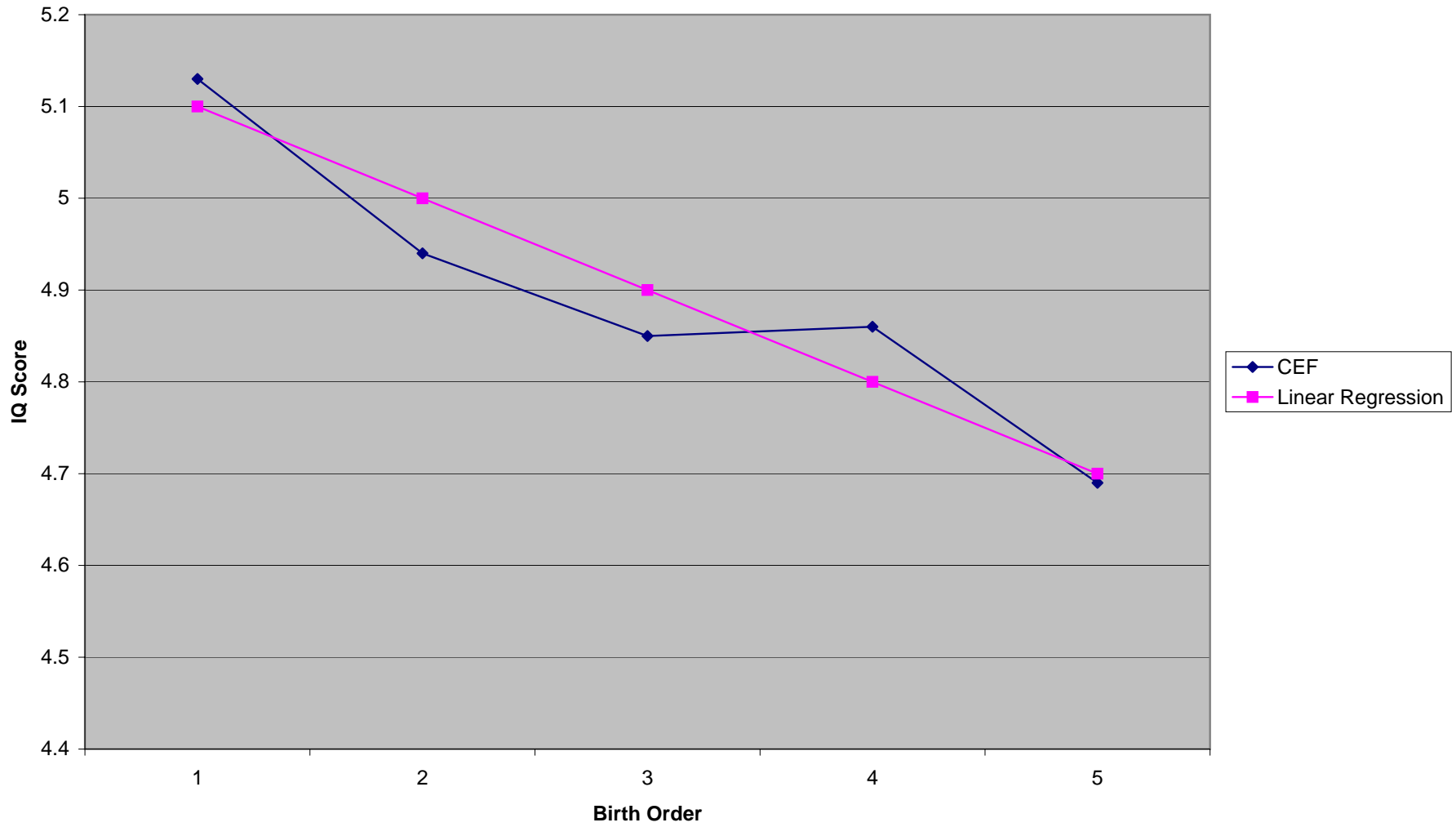
X_1	X_2	$E(Y/X)$	$X\hat{\beta}$	$E(\varepsilon/x)$
1	1	5.13	5.1	0.03
1	2	4.94	5	-0.06
1	3	4.85	4.9	-0.05
1	4	4.86	4.8	0.06
1	5	4.69	4.7	-0.01

2. X includes a constant and dummy variables for Birth Order = 2, 3, 4, 5.

$$\beta_1 = 5.13, \beta_2 = -0.19, \beta_3 = -0.28, \beta_4 = -0.27, \beta_5 = -0.44$$

X_1	X_2	X_3	X_4	X_5	$E(Y/X)$	$X\hat{\beta}$	$E(\varepsilon/x)$
1	0	0	0	0	5.13	5.13	0
1	1	0	0	0	4.94	4.94	0
1	0	1	0	0	4.85	4.85	0
1	0	0	1	0	4.86	4.86	0
1	0	0	0	1	4.69	4.69	0

Effect of Birth Order on IQ Score (5-child families)



Instrumental Variables

- The regression model is $Y_i = X_i'\beta + \varepsilon_i$
- OLS is consistent if $E(X_i\varepsilon_i) = 0$.
- If this assumption does not hold, OLS is *endogenous*.

Example 1: Simultaneous Equations

- The simplest supply and demand system:

$$q_i^s = \beta_s p_i + \varepsilon_s \quad (1)$$

$$q_i^d = \beta_d p_i + \varepsilon_d \quad (2)$$

$$q_i^s = q_i^d \quad (3)$$

- In equilibrium,

$$p_i = \frac{\varepsilon_d - \varepsilon_s}{\beta_s - \beta_d} \quad (4)$$

- Obviously, p_i is correlated with the error term in both equations.

Example 2: Omitted Variables

- Suppose $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$ but we exclude X_2 from the regression.
- Here X_1 is $n * K_1$ and X_2 is $n * (K - K_1)$.
- The new error term $v = X_2\beta_2 + \varepsilon$ is correlated with X_1 unless X_1 and X_2 are orthogonal or $\beta_2 = 0$.

Example 3: Measurement Error

- To see the effect of measurement error, consider the standard regression equation where there are no other control variables

$$y_i = \alpha + \beta x_i + \epsilon_i \quad (5)$$

- However, we observe

$$\tilde{x}_i = x_i + u_i \quad (6)$$

where u_i is mean zero and independent of all other variables. Substituting we get

$$y_i = \alpha + \beta(\tilde{x}_i - u_i) + \epsilon_i = \alpha + \beta\tilde{x}_i + v_i \quad (7)$$

- The new error term $v_i = \epsilon_i - \beta u_i$ is correlated with \tilde{x}_i .

Overview of Instrumental Variables

- The basics of IV can be understood with one x and one z .
- Consider the standard regression equation where there are no other control variables

$$y_i = \alpha + \beta x_i + \epsilon_i \quad (8)$$

- Let's define the sample covariance and variance matrices:

$$\text{cov}(x_i, y_i) = \frac{1}{n-1} \sum_i (x_i - \bar{x}_i)(y_i - \bar{y}_i) \quad (9)$$

$$\text{var}(x_i) = \frac{1}{n-1} \sum_i (x_i - \bar{x}_i)^2 \quad (10)$$

- The OLS estimator of β is

$$\hat{\beta}_{OLS} = \frac{\text{cov}(x_i, y_i)}{\text{var}(x_i)} = \frac{\text{cov}(x_i, \alpha + \beta x_i + \epsilon_i)}{\text{var}(x_i)} = \beta + \frac{\text{cov}(x_i, \epsilon_i)}{\text{var}(x_i)} \quad (11)$$

If x_i and ϵ_i are uncorrelated ($E(x_i\epsilon_i) = 0$), the probability limit of the second term is zero and OLS is a consistent estimator (as the sample size increases, the probability that the OLS estimate is not arbitrarily close to β goes to zero).

- However if x_i and ϵ_i are correlated ($E(x_i\epsilon_i) \neq 0$), the probability limit of $\text{cov}(x_i, \epsilon_i)$ does not equal zero and OLS is inconsistent.
- An *instrumental variable*, z_i , is one which is correlated with x_i but not with ϵ_i . The instrumental variables (IV) estimator of β is

$$\hat{\beta}_{IV} = \frac{\text{cov}(z_i, y_i)}{\text{cov}(z_i, x_i)} = \frac{\text{cov}(z_i, \alpha + \beta x_i + \epsilon_i)}{\text{cov}(z_i, x_i)} = \beta + \frac{\text{cov}(z_i, \epsilon_i)}{\text{cov}(z_i, x_i)} \quad (12)$$

- Given the assumption that z_i and ϵ_i are uncorrelated ($E(z_i\epsilon_i) = 0$), the probability limit of the second term is zero and the IV estimator is a consistent estimator.
- When there is only one instrument, the IV estimator can be calculated using the following procedure: (1) Regress x_i on z_i

$$x_i = \mu + \pi z_i + v_i \quad (13)$$

and form $\hat{x}_i = \hat{\mu} + \hat{\pi}z_i$

Then (2) estimate $\hat{\beta}$ by running the following regression by OLS

$$y_i = a + \beta \hat{x}_i + e_i \quad (14)$$

This process is called *Two Stage Least Squares (2SLS)*.

- It is quite easy to show this equivalence:

$$\hat{\beta}_{2SLS} = \frac{\text{cov}(\hat{x}_i, y_i)}{\text{var}(\hat{x}_i)} = \frac{\text{cov}(\hat{\mu} + \hat{\pi}z_i, y_i)}{\text{var}(\hat{\mu} + \hat{\pi}z_i)} \quad (15)$$

$$= \frac{\hat{\pi}\text{cov}(z_i, y_i)}{\hat{\pi}^2\text{var}(z_i)} = \frac{\text{cov}(z_i, y_i)}{\hat{\pi}\text{var}(z_i)} \quad (16)$$

Given

$$\hat{\pi} = \frac{\text{cov}(z_i, x_i)}{\text{var}(z_i)} \quad (17)$$

This implies that

$$\hat{\beta}_{2SLS} = \frac{\text{cov}(z_i, y_i)}{\text{cov}(z_i, x_i)} = \hat{\beta}_{IV} \quad (18)$$

- The *First Stage* refers to the regression $x_i = \mu + \pi z_i + v_i$.
- The *Reduced Form* refers to the regression $y_i = \theta + \delta z_i + u_i$.

- The *Indirect Least Squares (ILS)* estimator of β is $\hat{\delta}/\hat{\pi}$. This also equals $\hat{\beta}_{IV}$.

- To see this, note that

$$\hat{\beta}_{ILS} = \frac{\hat{\delta}}{\hat{\pi}} = \frac{\text{cov}(z_i, y_i)}{\text{var}(z_i)} \frac{\text{var}(z_i)}{\text{cov}(z_i, x_i)} = \frac{\text{cov}(z_i, y_i)}{\text{cov}(z_i, x_i)} \quad (19)$$

- OLS is often inconsistent because there are omitted variables. IV allows us to consistently estimate the coefficient of interest without actually having data on the omitted variables or even knowing what they are.
- Instrumental variables use only part of the variability in x – specifically, a part that is uncorrelated with the omitted variables – to estimate the relationship between x and y .

- A good instrument, z , is correlated with x for a clear reason, but uncorrelated with y for reasons beyond its effect on x .

Examples of Instruments

- Distance from home to nearest fast food as instrument for obesity.
- Twin births and sibling sex composition as instruments for family size
- Compulsory schooling laws as instruments for education.
- Tax rates on cigarettes as instruments for smoking.
- Weather shocks as instruments for income in developing countries
- Month of birth as instrument for school starting age.

The General Model

- With multiple instruments (overidentification), we could construct several IV estimators.
- 2SLS combines instruments to get a single more precise estimate.
- In this case, the instruments must all satisfy assumptions $E(z_i \epsilon_i) = 0$.
- We can write the models as

$$\begin{aligned} Y &= X\beta + \epsilon \\ X &= Z\pi + v. \end{aligned}$$

X is a matrix of exogenous and endogenous variables ($n \times K$).

Z is a matrix of exogenous variables and instruments ($n \times Q$), $Q \geq K$.

The 2SLS estimator is

$$\hat{\beta}_{2SLS} = (X'P_ZX)^{-1}X'P_ZY \quad (20)$$

where

$$P_Z = Z(Z'Z)^{-1}Z' \quad (21)$$

- It can be shown that the 2SLS estimator is the most efficient IV estimator.
- The **Order Condition** for identification is that there must be at least as many instruments as endogenous variables: $Q \geq K$. This is a necessary but not sufficient condition.
- The **Rank Condition** for identification is that $rank(Z'X) = K$. This is a sufficient condition and it ensures that there is a first stage relationship.

- Example: Suppose we have 2 endogenous variables

$$\begin{aligned}x_1 &= a_1 + a_2z_1 + a_3z_2 + u_1 \\x_2 &= b_1 + b_2z_1 + b_3z_2 + u_2\end{aligned}$$

The order condition is satisfied. However, if $a_2 = 0$ and $b_2 = 0$, the rank condition fails and the model is unidentified. If $a_2 = 0$ and $b_3 = 0$ and the other parameters are non-zero, the rank condition passes and the model is identified.

Variance of 2SLS Estimator

- Recall the 2SLS estimator

$$\begin{aligned}\hat{\beta}_{2SLS} &= (X'P_ZX)^{-1}X'P_Zy \\ &= (\hat{X}'\hat{X})^{-1}\hat{X}'Y\end{aligned}$$

where $\hat{X} = P_ZX$ is the predicted value of X from the first stage regression.

- Given this is just the parameter from an OLS regression of Y on \widehat{X} , the estimated covariance matrix under homoskedasticity takes the same form as OLS:

$$\widehat{Var}(\widehat{\beta}_{2SLS}) = \widehat{\sigma}^2 (\widehat{X}'\widehat{X})^{-1} = \widehat{\sigma}^2 (X'P_Z X)^{-1}$$

where

$$\widehat{\sigma}^2 = \frac{1}{n - K} (Y - X\widehat{\beta}_{2SLS})'(Y - X\widehat{\beta}_{2SLS})$$

- Note that $\widehat{\sigma}^2$ uses X rather than \widehat{X} . Shows that standard errors from doing 2SLS manually are incorrect.
- We can simplify this further in the case of the bivariate model from equations (8) and (13).

- In this case, the element in the second row and second column of $X'P_ZX = (X'Z)(Z'Z)^{-1}(Z'X)$ simplifies to (algebra is a bit messy)

$$\frac{n^2 \text{cov}(z_i, x_i)^2}{n \text{var}(z_i)}$$

implying that the relevant element of

$$(X'P_ZX)^{-1} = \frac{1}{n\rho_{xz}^2\sigma_x^2} \quad (22)$$

where the correlation between x and z equals

$$\rho_{xz} = \frac{\text{cov}(z_i, x_i)}{\sigma_x\sigma_z}$$

- Equation (22) tells us that the 2SLS variance

1. Decreases at a rate of $1/n$.

2. Decreases as the variance of the explanatory variable increases.
3. Decreases with the correlation between x and z . If this correlation approaches zero, the 2SLS variance goes to infinity.
4. Is higher than the OLS variance as, for OLS, $\rho_{xz} = 1$ as OLS uses x as an instrument for itself.

Hausman Tests

- Also referred to as Wu-Hausman, or Durbin-Wu-Hausman tests.
- Have wide applicability to cases where there are two estimators and

1. Estimator 1 is consistent and efficient under the null but inconsistent under the alternative.

2. Estimator 2 is consistent in either case but is inefficient under the null.
 - We will only consider 2SLS and OLS cases.

 - The null hypothesis is that $E(X'\varepsilon) = 0$.

 - Suppose we have our model $Y = X\beta + \varepsilon$.

 - If $E(X'\varepsilon) = 0$ the OLS estimator provides consistent estimates.

- If $E(X'\varepsilon) \neq 0$ and we have valid instruments, 2SLS is consistent but OLS is not.
- If $E(X'\varepsilon) = 0$ 2SLS remains consistent but is less efficient than OLS.

- Hausman suggests the following test statistic for whether OLS is consistent:

$$h = (\hat{\beta}_{OLS} - \hat{\beta}_{2SLS})' [V(\hat{\beta}_{2SLS}) - V(\hat{\beta}_{OLS})]^{-1} (\hat{\beta}_{OLS} - \hat{\beta}_{2SLS})$$

which has an asymptotic chi square distribution.

- Note that a nice feature is that one does not need to estimate the covariance of the two estimators.

Hausman Test as Vector of Contrasts (1)

- Compare the OLS estimator $\hat{\beta}_{OLS} = (X'X)^{-1}X'Y$ to the 2SLS estimator $\hat{\beta}_{2SLS} = (X'P_ZX)^{-1}X'P_ZY$ where P_Z is symmetric $n * n$ matrix with rank of at least K .
- Under the null hypothesis $E(X'\varepsilon) = 0$, both are consistent.

$$\begin{aligned}
 \hat{\beta}_{2SLS} - \hat{\beta}_{OLS} &= (X'P_ZX)^{-1}X'P_ZY - (X'X)^{-1}X'Y \\
 &= (X'P_ZX)^{-1} \left[X'P_ZY - (X'P_ZX)(X'X)^{-1}X'Y \right] \\
 &= (X'P_ZX)^{-1}X'P_Z \left[I - X(X'X)^{-1}X' \right] Y \\
 &= (X'P_ZX)^{-1}X'P_ZM_XY \tag{23}
 \end{aligned}$$

- The probability limit of this difference will be zero when

$$p \lim \frac{1}{n} X'P_ZM_XY = 0 \tag{24}$$

- We can partition the X matrix as $X = [X_1 X_2]$ where X_1 is an $n * G$ matrix of potentially endogenous variables and X_2 is an $n * (K - G)$ matrix of exogenous variables.
- We have instruments Z where $Z = [Z^* X_2]$ an $n * Q$ matrix ($Q \geq K$).
- Letting hats denote the first stage predicted values, clearly $\widehat{X}_2 = X_2$ and $X_2 M_x$ is zero for the rows of M_x corresponding to X_2 .
- Therefore checking that $p \lim \frac{1}{n} X' P_z M_X Y = 0$ reduces to checking whether $p \lim \frac{1}{n} X_1' P_z M_X Y = p \lim \frac{1}{n} \widehat{X}_1' M_X Y = 0$.
- We can implement this test using an F-test on δ in the regression:

$$Y = X\beta + \widehat{X}_1\delta + error \quad (25)$$

- Denoting $\delta = 0$ as the restricted model, the F-statistic is

$$H = \frac{RSS_r - RSS_u}{RSS_u / (n - K - G)} \quad (26)$$

- Note from (23), that we can also do the test by regressing $M_X Y$ on \widehat{X} and testing whether the parameters are zero.

Hausman Test as Vector of Contrasts (2)

- Compare the OLS estimator $\widehat{\beta}_{OLS} = (X'X)^{-1}X'Y$ to a different OLS estimator where Z^* is added as a control:

$$Y = X\beta + Z^*\gamma + v \quad (27)$$

- Because of the exclusion restriction, Z^* should have no explanatory power when X is exogenous.

- Using the Frisch-Waugh-Lovell theorem, the resulting estimate of β is $\hat{\beta}_p$:

$$\hat{\beta}_p = (X' M_{Z^*} X)^{-1} X' M_{z^*} Y \quad (28)$$

- Subtracting the OLS estimator,

$$\begin{aligned} \hat{\beta}_p - \hat{\beta}_{OLS} &= (X' M_{z^*} X)^{-1} X' M_{z^*} Y - (X' X)^{-1} X' Y \\ &= (X' M_{z^*} X)^{-1} \left[X' M_{z^*} Y - (X' M_{z^*} X) (X' X)^{-1} X' Y \right] \\ &= (X' M_{z^*} X)^{-1} X' M_{z^*} \left[I - X (X' X)^{-1} X' \right] Y \\ &= (X' M_{z^*} X)^{-1} X' M_{z^*} M_X Y \end{aligned} \quad (29)$$

- By analogy with equation (25), we can implement an F-test for whether $p \lim \frac{1}{n} X' M_{z^*} M_X Y = 0$ by testing whether $\delta = 0$ in the regression:

$$Y = X\beta + M_{z^*} X\delta + \text{error} \quad (30)$$

- One can show (not easily) that the resulting F-statistic is identical to that above.
- Alternatively, we can regress $M_X Y$ on $M_{z^*} X$ and test whether the parameters are zero.

Two-Sample 2SLS

- Suppose have 2 samples from the same population.
- X is a matrix of exogenous and endogenous variables ($N \times K$). Z is a matrix of exogenous variables and instruments ($N \times Q$), $Q \geq K$.
- Sample 1 contains Y and Z but not the endogenous elements of X .
- Sample 2 contains X and Z but not Y .
- Using subscript to denote sample, can implement TS2SLS by

1. Do the first stage regression(s) using Sample 2 observations

$$X_2 = Z_2\pi + v_2 \quad (31)$$

and estimate

$$\hat{\pi} = (Z_2'Z_2)^{-1}Z_2'X_2 \quad (32)$$

2. Form the predicted value of X_1 as

$$\hat{X}_1 = Z_1\hat{\pi} \quad (33)$$

3. Regress Y on \hat{X}_1 using Sample 1 observations.

$$\hat{\beta}_{TS2SLS} = (\hat{X}_1'\hat{X}_1)^{-1}\hat{X}_1'Y \quad (34)$$

where

$$\hat{X}_1 = Z_1(Z_2'Z_2)^{-1}Z_2'X_2 \quad (35)$$

- Note that with one endogenous variable and one instrument, we can take an ILS approach.

1. Do step 1. as with TS2SLS to get $\hat{\pi}$.

2. Estimate the Reduced Form using Sample 1

$$Y = Z_1\delta + u_1 \quad (36)$$

3. Take the ratio of the Reduced Form and First Stage Estimates:

$$\hat{\beta}_{ILS} = \frac{\hat{\delta}}{\hat{\pi}} \quad (37)$$

- To see this note that

$$\begin{aligned}
 \hat{\beta}_{TS2SLS} &= (\widehat{X}'_1 \widehat{X}_1)^{-1} \widehat{X}'_1 Y \\
 &= [X'_2 Z_2 (Z'_2 Z_2)^{-1} Z'_1 Z_1 (Z'_2 Z_2)^{-1} Z'_2 X_2]^{-1} X'_2 Z_2 (Z'_2 Z_2)^{-1} Z'_1 Y \\
 &= (Z'_2 X_2)^{-1} (Z'_2 Z_2) (Z'_1 Z_1)^{-1} (Z'_2 Z_2) (X'_2 Z_2)^{-1} X'_2 Z_2 (Z'_2 Z_2)^{-1} Z'_1 Y \\
 &= (Z'_2 X_2)^{-1} (Z'_2 Z_2) (Z'_1 Z_1)^{-1} Z'_1 Y \\
 &= [(Z'_2 Z_2)^{-1} Z'_2 X_2]^{-1} (Z'_1 Z_1)^{-1} Z'_1 Y \\
 &= \hat{\pi}^{-1} \hat{\delta}
 \end{aligned} \tag{3}$$

- This derivation is valid so long as $Q = K$ so $X'_2 Z_2$ and $Z'_2 X_2$ are square matrices that we can invert.
- It also uses the matrix inversion rule: $(ABC)^{-1} = C^{-1} B^{-1} A^{-1}$

Example: Devereux and Hart 2010

- Return to education using UK change in compulsory school law.
- If born 1933 or after have minimum leaving age of 15 rather than 14.
- New Earnings Survey has earnings and cohort (Y and Z).
- General Household Survey has education and cohort (X and Z).
- In General Household Survey estimate

$$Education = \alpha_0 + \alpha_1 \mathbf{1}(YOB \geq 1933) + W\alpha_2 + e_1 \quad (39)$$

- In New Earnings Survey estimate

$$\text{Log}(\text{earnings}) = \gamma_0 + \gamma_1 \mathbf{1}(YOB \geq 1933) + W\gamma_2 + e_2 \quad (40)$$

- Then the TS2SLS estimator of the return to education is $\hat{\beta} = \frac{\hat{\gamma}_1}{\alpha_1}$.
- To calculate the standard error we use the delta method.

The Delta Method

- This is a method for estimating variances of functions of random variables using Taylor-series expansions.

$$f(x, y) = f(x_0, y_0) + \frac{\partial f(x, y)}{\partial x} \Big|_{x_0, y_0} (x - x_0) + \frac{\partial f(x, y)}{\partial y} \Big|_{x_0, y_0} (y - y_0) + \dots$$

- For the case where $f(x, y) = y/x$, $\frac{\partial f(x, y)}{\partial x} = \frac{-y}{x^2}$ and $\frac{\partial f(x, y)}{\partial y} = \frac{1}{x}$.
- Therefore, evaluating at the means of x and y ,

$$\frac{y}{x} \simeq \frac{\mu_y}{\mu_x} - \frac{\mu_y}{\mu_x^2}(x - \mu_x) + \frac{1}{\mu_x}(y - \mu_y) \quad (41)$$

- Then,

$$\text{var}\left(\frac{y}{x}\right) \simeq \frac{\mu_y^2}{\mu_x^4} \text{var}(x) + \frac{1}{\mu_x^2} \text{var}(y) - 2 \frac{\mu_y}{\mu_x^3} \text{cov}(x, y) \quad (42)$$

- In our case

$$\text{var}(\hat{\beta}) \simeq \frac{\hat{\gamma}_1^2}{\hat{\alpha}_1^4} \text{var}(\hat{\alpha}_1) + \frac{1}{\hat{\alpha}_1^2} \text{var}(\hat{\gamma}_1) \quad (43)$$

- Note that the covariance term disappears because the parameters are estimated from 2 independent samples.

The Method of Moments (MOM)

- A *population moment* is just the expectation of some continuous function of a random variable:

$$\gamma = E[g(x_i)] \quad (1)$$

- For example, one moment is the mean: $\mu = E(x_i)$.
- The variance is a function of two moments:

$$\sigma^2 = E[x_i - E(x_i)]^2 \quad (2)$$

$$= E(x_i^2) - [E(x_i)]^2 \quad (3)$$

- We also refer to functions of moments as moments.

- A *sample moment* is the analog of a population moment from a particular random sample

$$\hat{\gamma} = \frac{1}{n} \sum_i g(x_i) \quad (4)$$

- So, the sample mean is $\hat{\mu} = \frac{1}{n} \sum_i x_i$.
- The idea of MOM is to estimate a population moment using the corresponding sample moment.
- For example, the MOM estimator of the variance using (3) is

$$\hat{\sigma}^2 = \left(\frac{1}{n} \sum_i x_i^2 \right) - \left[\frac{1}{n} \sum_i x_i \right]^2 \quad (5)$$

$$= \frac{1}{n} \sum_i (x_i - \bar{x}_i)^2 \quad (6)$$

- This is very similar to our usual estimator of the variance

$$\frac{1}{n-1} \sum_i (x_i - \bar{x}_i)^2 \quad (7)$$

- The MOM estimator is biased but is consistent.
- Alternatively, we could calculate the MOM estimator directly using (2)

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \bar{x}_i)^2 \quad (8)$$

OLS as Methods of Moments Estimator

- Our population parameters for linear regression were

$$\beta = E(X_i X_i')^{-1} E(X_i Y_i) \quad (9)$$

- Can derive method of moments estimator by replacing population moments $E(X_i X_i')$ and $E(X_i Y_i)$ by sample moments:

$$b = \left[\frac{1}{n} \sum_i X_i X_i' \right]^{-1} \frac{1}{n} \sum_i X_i Y_i = (X'X)^{-1} X'Y \quad (10)$$

- Or alternatively, we can use the population moment condition

$$E(X_i \varepsilon_i) = E(X_i (Y_i - X_i' \beta)) \quad (11)$$

- The MOM approach is to choose an estimator b so that it sets the sample analog of (11) to zero:

$$\frac{1}{n} \sum_i X_i (Y_i - X_i' b) = 0 \quad (12)$$

This implies that

$$\frac{1}{n} \sum_i X_i Y_i = \frac{1}{n} \sum_i X_i X_i' b \quad (13)$$

So

$$b = \left[\frac{1}{n} \sum_i X_i X_i' \right]^{-1} \frac{1}{n} \sum_i X_i Y_i = (X'X)^{-1} X'Y \quad (14)$$

- Note that this is the OLS estimator.

Generalized Method of Moments (GMM)

- We saw earlier that the OLS estimator solves the moment condition

$$E(X_i(Y_i - X_i'\beta)) = 0 \quad (15)$$

- This moment condition was motivated by the condition $E(X_i\varepsilon_i) = 0$.
- This type of approach can be extended.
- For example, we may know that $E(Z_i\varepsilon_i) = 0$ where Z_i may include some of the elements of X_i .
- The idea of GMM is to substitute out the error term with a function of data and parameters.
- Then find the parameter values that make the conditions hold in the sample.

- Let $\varepsilon_i(\beta) = (Y_i - X_i'\beta)$. We find the parameter such that

$$\frac{1}{n} \sum_i g_i(\beta) = \frac{1}{n} \sum_i Z_i \varepsilon_i(\beta) = \frac{1}{n} Z'(Y - X\beta) \quad (16)$$

is as close as possible to zero.

- A first guess might be the MOM estimator

$$\hat{\beta} = (Z'X)^{-1}Z'Y \quad (17)$$

but this only works if $Z'X$ is invertible and this is only the case if it is a square matrix.

- MOM only works when the number of moment conditions equals the number of parameters to be estimated.

- Instead GMM solves the following problem:

$$\min \left(\frac{1}{n} Z'(Y - X\beta) \right)' W \left(\frac{1}{n} Z'(Y - X\beta) \right) \quad (18)$$

- Here W is called the weight matrix and is some positive definite (PD) square matrix.
- Taking the first order conditions, we get

$$\hat{\beta} = (X'ZWZ'X)^{-1}X'ZWZ'Y \quad (19)$$

- To see this, note that $(Z'(Y - X\beta))' W (Z'(Y - X\beta))$

$$\begin{aligned} &= (Z'Y - Z'X\beta)'W(Z'Y - Z'X\beta) \\ &= Y'ZWZ'Y - Y'ZWZ'X\beta - \beta'X'ZWZ'Y + \beta'X'ZWZ'X\beta \\ &= Y'ZWZ'Y - 2\beta'X'ZWZ'Y + \beta'X'ZWZ'X\beta \end{aligned}$$

This uses the fact that the transpose of a scalar is itself. Then, taking first order conditions

$$-2X'ZWZ'Y + 2X'ZWZ'X\beta = 0 \quad (20)$$

- $X'ZWZ'X$ will be invertible so long as the number of moment conditions, Q (elements of Z) is as least as big as the number of parameters, K (elements of X).
- For example, not invertible if

$$Y_i = X_{1i}\beta_1 + X_{2i}\beta_2 + \varepsilon_i \quad (21)$$

$$Z_i = X_{1i}\alpha_1 \quad (22)$$

- When $Q > K$, GMM estimates will not cause all moment conditions to equal zero but will get them as close to zero as possible.

- When $Q = K$, as we would expect from (17),

$$\hat{\beta} = (Z'X)^{-1}Z'Y \quad (23)$$

To see this note that when $Q = K$, $Z'X$ is a square matrix so

$$(X'ZWZ'X)^{-1} = (Z'X)^{-1}W^{-1}(X'Z)^{-1} \quad (24)$$

(remember $(ABC)^{-1} = C^{-1}B^{-1}A^{-1}$). Also note that in this case, W plays no role.

- This is exactly the IV estimator we saw earlier.
- If $X = Z$, the GMM estimator is exactly the OLS estimator.

Consistency of GMM Estimator

- The GMM estimator

$$\hat{\beta} = (X'ZWZ'X)^{-1}X'ZWZ'Y \quad (25)$$

$$= \beta + (X'ZWZ'X)^{-1}X'ZWZ'\varepsilon \quad (26)$$

$$= \beta + \left(\frac{X'Z}{n}W\frac{Z'X}{n}\right)^{-1} \left(\frac{X'Z}{n}W\frac{Z'\varepsilon}{n}\right) \quad (27)$$

- Using the Law of Large Numbers (LLN),

$$\frac{X'Z}{n} = 1/n \sum_i X_i Z_i' \rightarrow \Sigma_{XZ} \quad (28)$$

$$\frac{Z'X}{n} = 1/n \sum_i Z_i X_i' \rightarrow \Sigma_{ZX} \quad (29)$$

$$\frac{Z'\varepsilon}{n} = 1/n \sum_i Z_i \varepsilon_i \rightarrow E(Z_i \varepsilon_i) \quad (30)$$

- Denote

$$H = (\Sigma_{XZ}W\Sigma_{ZX})^{-1} \Sigma_{ZX}W \quad (31)$$

- Then

$$\hat{\beta} - \beta \rightarrow HE(Z_i\varepsilon_i) = 0 \quad (32)$$

showing consistency of GMM for any PD weighting matrix, W .

Choice of Weight Matrix

- Under some regulatory conditions, the GMM $\hat{\beta}$ is also asymptotically normally distributed for any PD W .
- If the model is overidentified ($Q > K$), the choice of weight matrix affects the asymptotic variance and also the coefficient estimates in finite samples.

- The "best" choice for W is the inverse of the covariance of the moments i.e. the inverse of the covariance matrix of

$$Z'(Y - X\beta) = \sum_i Z_i \varepsilon_i \quad (33)$$

- However, this is unknown and needs to be estimated in the data. We can use a 3-step procedure

1. Choose a weight matrix and do GMM. Any PD weighting matrix will give consistent estimates. A good initial choice is

$$W = (Z'Z/n)^{-1} \quad (34)$$

This gives the estimator

$$\hat{\beta} = (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \quad (35)$$

$$= (X'P_zX)^{-1}X'P_zY \quad (36)$$

This is exactly the 2SLS estimator we saw earlier.

2. Take the residuals and use them to estimate the covariance of the moments

$$\widehat{Var}\left(\sum_i Z_i \varepsilon_i\right) = \frac{1}{n} \sum_i e_i^2 Z_i Z_i' \quad (37)$$

where

$$e_i = Y_i - X_i' \hat{\beta} \quad (38)$$

3. Do GMM with \widehat{W} as weight matrix where $\widehat{W} = \left(\frac{1}{n} \sum_i e_i^2 Z_i Z_i'\right)^{-1}$. Low variance moments are given higher weight in estimation than high variance moments.

- Note that if the errors are homoskedastic, this is just the 2SLS estimator.

Variance of GMM Estimator

- In the general case, the GMM estimator is

$$\min g(\beta)'Wg(\beta) \quad (39)$$

where the moment conditions are $g(\beta) = 0$.

- The variance of the GMM estimator is

$$\frac{1}{n}(G'WG)^{-1}G'W\Psi WG(G'WG)^{-1} \quad (40)$$

where $G = \frac{\partial g(\beta)}{\partial \beta}$ and Ψ is the variance-covariance matrix of the moments.

- In the OLS case,

$$g(\beta) = E(X_i \varepsilon_i) = 0 \quad (41)$$

$$\hat{g}(\beta) = 1/n \sum_i X_i (Y_i - X_i' \beta) = 0 \quad (42)$$

$$W = I \quad (43)$$

$$G = \frac{\partial \hat{g}(\beta)}{\partial \beta} = \frac{X' X}{n} \quad (44)$$

$$\Psi = E[(X_i \varepsilon_i)(X_i \varepsilon_i)'] = E[\varepsilon_i^2 X' X] \quad (45)$$

$$\hat{\Psi} = \hat{\sigma}_\varepsilon^2 \frac{X' X}{n} \quad (46)$$

where the last step assumes homoskedasticity. Putting these together we get

$$\hat{V}(\hat{\beta}) = \hat{\sigma}_\varepsilon^2 (X' X)^{-1} \quad (47)$$

- In the 2SLS case,

$$g(\beta) = E(Z_i \varepsilon_i) = 0 \quad (48)$$

$$\hat{g}(\beta) = \frac{1}{n} \sum_i Z_i (Y_i - X_i' \beta) = 0 \quad (49)$$

$$W = \left(\frac{Z'Z}{n} \right)^{-1} \quad (50)$$

$$G = \frac{\partial \hat{g}(\beta)}{\partial \beta} = \frac{Z'X}{n} \quad (51)$$

$$\Psi = E[(Z_i \varepsilon_i)(Z_i \varepsilon_i)'] = E[\varepsilon_i^2 Z'Z] \quad (52)$$

$$\hat{\Psi} = \hat{\sigma}_\varepsilon^2 \frac{Z'Z}{n} \quad (53)$$

where the third and last steps assume homoskedasticity. Putting these together we get

$$\hat{V}(\hat{\beta}) = \hat{\sigma}_\varepsilon^2 (X'P_Z X)^{-1} \quad (54)$$

- This formula ignores the fact that the weight matrix is estimated and so

may understate the true variance.

Why Use GMM in Linear Models?

- When the model is just identified, GMM coincides with IV or OLS. So no reason to use GMM.

- In overidentified models with homoskedastic errors,

$$\widehat{W} = \left(\frac{1}{n} \sum_i e_i^2 Z_i Z_i' \right)^{-1} = \hat{\sigma}_e^2 \frac{1}{n} \sum_i (Z_i Z_i')^{-1}$$

and the GMM estimator coincides with 2SLS. So no reason to use GMM.

- In overidentified models with heteroskedasticity, GMM is more efficient than 2SLS.

- Also, in time series models with serial correlation, GMM is more efficient than 2SLS.
- When estimating a system of equations, GMM is particularly useful. You will see this in Kevin Denny's section of the course.

Relationship of GMM to Maximum Likelihood

Maximum Likelihood Interpretation of OLS

- The regression model is

$$Y_i = X_i' \beta + \varepsilon_i \quad (55)$$

- Assume that

$$\varepsilon_i / X_i \sim i.i.d. N(0, \sigma^2) \quad (56)$$

$$X_i \sim i.i.d. g(x) \quad (57)$$

- The likelihood function is the joint density of the observed data evaluated at the observed data values.
- The joint density of $\{Y_i, X_i\}$ is

$$f(y, x) = f(y/x)g(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(y - x'\beta)^2\right\} g(x) \quad (58)$$

- The likelihood function is

$$L = \prod \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(Y_i - X_i'\beta)^2\right\} g(X_i) \quad (59)$$

$$= \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left\{-\frac{1}{2\sigma^2} \sum_i (Y_i - X_i'\beta)^2\right\} \prod g(X_i) \quad (60)$$

- Taking Logs,

$$\text{Log}L = -\frac{n}{2}\text{Log}(\sigma^2) - \frac{n}{2}\text{Log}(2\pi) - \frac{1}{2\sigma^2} \sum_i (Y_i - X_i'\beta)^2 + \sum_i \text{Log}g(X_i)$$

- Ignoring the last term which is not a function of β or σ ,

$$\begin{aligned} \text{Log}L &= -\frac{n}{2}\text{Log}(\sigma^2) - \frac{n}{2}\text{Log}(2\pi) - \frac{1}{2\sigma^2}(Y - X\beta)'(Y - X\beta) \\ &= -\frac{n}{2}\text{Log}(\sigma^2) - \frac{n}{2}\text{Log}(2\pi) - \frac{1}{2\sigma^2} \{Y'Y - 2\beta'X'Y + \beta'X'X\beta\} \end{aligned}$$

- Taking first order conditions of this scalar with respect to β and σ^2

$$\frac{1}{\hat{\sigma}^2} \{X'Y + X'X\hat{\beta}\} = 0 \quad (61)$$

$$-\frac{n}{2\hat{\sigma}^2} + \frac{1}{2\hat{\sigma}^4}(Y - X\hat{\beta})'(Y - X\hat{\beta}) = 0 \quad (62)$$

- These imply the MLE

$$\hat{\beta} = (X'X)^{-1}X'Y \quad (63)$$

$$\hat{\sigma}^2 = \frac{(Y - X\hat{\beta})'(Y - X\hat{\beta})}{n} \quad (64)$$

- Note that when taking the first FOC, we are just minimising the sum of squared errors.

GMM Interpretation of Maximum Likelihood

- In the general case, the GMM estimator is

$$\min g(\beta)'Wg(\beta) \quad (65)$$

where the moment conditions are $g(\beta) = 0$.

- The FOC are

$$2W \frac{\partial g}{\partial \beta} g(\beta) = 0 \quad (66)$$

- Consider the following moment:

$$g(\beta) = \frac{\partial \text{Log} L}{\partial \beta} \quad (67)$$

so

$$\frac{\partial g}{\partial \beta} = \frac{\partial^2 \text{Log} L}{\partial \beta \partial \beta'} \quad (68)$$

- The optimal weight matrix is the inverse of the variance-covariance matrix of the moments. In this case,

$$V[g(\beta)] = V \left[\frac{\partial \text{Log} L}{\partial \beta} \right] = -E \left[\frac{\partial^2 \text{Log} L}{\partial \beta \partial \beta'} \right] \quad (69)$$

and, so, the best estimate of the optimal weighting matrix is

$$\left(\frac{\partial^2 \text{Log}L}{\partial \beta \partial \beta'} \right)^{-1} \quad (70)$$

- Substituting these into the FOC, we find that the GMM estimator is defined by $\frac{\partial \text{Log}L}{\partial \beta} = 0$, the same as ML.
- So the ML estimator can be seen as a GMM estimator with a particular set of moment equations.

Limited Information Maximum Likelihood (LIML)

- ML version of 2SLS.

- Assumes joint normality of the error terms.
- LIML estimate exactly the same as 2SLS if model is just identified.
- LIML and 2SLS are asymptotically equivalent.
- LIML has better small sample properties than 2SLS in over-identified models.