

ECON 6200
Econometrics Notes

Gabe Sekeres

Spring 2025

Contents

1	The Linear Model	2
2	The Linear Model in Large Samples	12
3	Instrumental Variables	16
4	Generalized Method of Moments	21
5	Panel Data	31
6	Extremum Estimation	33
7	Worked Examples	45
8	Bootstrapping	52
9	Non-Parametrics	54

Introduction

This course builds directly on ECON 6190, [the material](#) will not be revised. If you need to catch up on the material, especially probability and statistics, see the [Hansen](#) and [Durrett](#) textbooks. If there's one econometrics book to actually own, it's the Hansen. These notes will often reference results in either Hansen or [Hayashi](#).

See the syllabus for how we will be assessed – homework is 30%, prelim is 30%, and the final is 40%. For homework, you are invited (and *highly* recommended) to form study groups! It's completely okay to work on the homeworks in study groups, but you will submit individual write-ups.

Roughly, in the first few weeks, we will start with the linear model in some detail. We will then generalize to IV, TSLS, and go rapidly to the Generalized Method of Moments (GMM) and extremum (or m -) estimation. We will also cover some nonparametrics as well as bootstrap.

You can think of us starting in the most specific case and moving outwards. OLS is a special case of IV, which is a special case of TSLS, which is a special case of GMM, which is a special case of extremum estimation. We will also think about Maximum Likelihood (ML) as a special case of extremum estimation, and panel data as a special case of GMM. Studying this tree will cover half of the course, and then we will have some time to cover nonparametric methods like kernel density and kernel mean, as well as bootstrapping. This may seem like we are doing some old school stuff, since the other ML (machine learning) has overtaken these methods, but theoretically they are very similar.

1 The Linear Model

We have, of course, already encountered the

Definition. *Ordinary Least Squares (OLS)* estimator:

$$\begin{aligned}\hat{\beta} &= (X'X)^{-1} X'Y && \text{Data Matrix} \\ &= (\mathbb{E}_n XX')^{-1} \mathbb{E}_n XY && \text{Sample Expectation} \\ &= \left(\frac{1}{n} \sum_{i=1}^n X_i X_i' \right) \frac{1}{n} \sum_{i=1}^n X_i Y_i && \text{Sample Average} \\ &= \left(\sum_{i=1}^n X_i X_i' \right) \sum_{i=1}^n X_i Y_i && \text{Clean Sample Average}\end{aligned}$$

All of these notations are equivalent, but helpful in different concepts. Data matrix notation is extremely useful for proving results. The sample expectation is useful because it reminds us of the empirical context – this is a shorthand for sample average, which is written in two ways depending on context – essentially, when we care about the asymptotic behavior in different ways.

The interpretation of $\hat{\beta}$ depends on context:

1. In any given sample, it just projects Y onto X
2. Under weak assumptions, it converges to the population analog $\beta^* \equiv (\mathbb{E}XX')^{-1}\mathbb{E}XY$, which is the population projection coefficient and characterizes the *best linear predictor under square loss*
3. Under stronger assumptions, it estimates a causal effect of X on Y .

We will elaborate these in order, and develop the classic theory of Least Squares estimation.

Recall that $\hat{\beta}$ can be derived as the minimization (in b) of

$$\sum_{i=1}^n (Y_i - X_i' b)^2 = (Y - Xb)'(Y - Xb)$$

which is of course where the name comes from. However, recall that also this minimization defined b such that Xb is the point in the span of X that is *closest* to Y in Euclidean distance. Basically, we projected Y onto X .

We can see this with the illustration in Figure 1, which uses demeaned vectors, and defines the projection $\hat{Y} \equiv X\hat{\beta} = \beta_1 X_1 + \beta_2 X_2$ and the residual $\hat{\varepsilon} \equiv Y - \hat{Y}$.

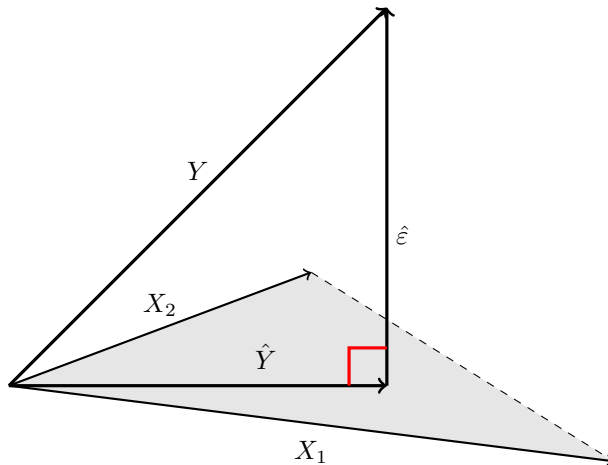


Figure 1: Illustration of OLS as Projection

Corollary 1.1. Sum of Squares Decomposition *It immediately follows from Pythagoras that*

$$Y'Y = \hat{Y}'\hat{Y} + \hat{\varepsilon}'\hat{\varepsilon}$$

Equivalently, $SST = SSE + SSR$ or

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{\varepsilon}_i^2$$

It follows immediately that $R^2 = SSE / SST \in [0, 1]$. The extreme values of R^2 correspond to $\hat{Y} = 0$ and $\hat{\varepsilon} = 0$ respectively. It also follows that $\hat{\varepsilon}$ is by construction orthogonal to \hat{Y} .

We will recap this and basic related facts in the first homework.

Question. What happens with collinear X ?

With collinear X , the set on which we project is of lower dimension. The projection \hat{Y} is still unique, but the projection coefficient is not.

Definition. The *projection coefficient* $\hat{\beta}$ is defined as

$$\hat{\beta} \equiv \operatorname{argmin}_b \sum_i (Y_i - X_i' b)^2 = \operatorname{argmin}_b (Y - Xb)'(Y - Xb)$$

and can be characterized by the first order condition

$$\begin{aligned}\frac{\partial}{\partial b}(Y - Xb)'(Y - Xb) &= -2X'Y + 2X'Xb \stackrel{!}{=} 0 \\ \implies \hat{\beta} &= (X'X)^{-1}X'Y\end{aligned}$$

The fitted values and residuals equal

$$\begin{aligned}\hat{Y} = X\hat{\beta} &= \underbrace{X(X'X)^{-1}X'}_{P_X, \text{ the projection matrix}} Y = P_X Y \\ \hat{\varepsilon} = Y - \hat{Y} &= \underbrace{(I_n - X(X'X)^{-1}X')}_{\text{annihilator matrix}} Y\end{aligned}$$

Example. Frisch-Waugh(-Lovell) The projection of Y on X can be decomposed, giving rise to some important results. To fix ideas, consider projecting Y :

- on the scalar X_1 (plus a constant), getting slope coefficient $\tilde{\beta}_1$, versus
- on the scalars (X_1, X_2) (plus a constant), getting slope coefficients $(\hat{\beta}_1, \hat{\beta}_2)$.

Can we interestingly relate $\tilde{\beta}_1$ to $\hat{\beta}_1$? Yes! To do so, consider projecting X_1 on X_2 , getting slope coefficient $\hat{\gamma}$. To simplify, assume all variables are demeaned. Across regressions, this leads to the first order conditions:

$$\begin{aligned}0 &= \mathbb{E}_n \left(X_1 \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_2 X_2 \right) \right) \\ 0 &= \mathbb{E}_n \left(X_2 \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_2 X_2 \right) \right) \\ 0 &= \mathbb{E}_n \left(X_2 \left(X_1 - \hat{\gamma} X_2 \right) \right)\end{aligned}$$

We can combine the first two, and use the third to get

$$\begin{aligned}0 &= \mathbb{E}_n(X_1 - \hat{\gamma}X_2) \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_2 X_2 \right) \\ &= \mathbb{E}_n(X_1 - \hat{\gamma}X_2) \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_1 \hat{\gamma} X_2 \right) \\ &= \mathbb{E}_n(X_1 - \hat{\gamma}X_2) \left(Y - \hat{\beta}_1 (X_1 - \hat{\gamma}X_2) \right)\end{aligned}$$

The last line is the first order condition from regressing Y on the residuals of the regression of X_1 on X_2 . Thus, $\hat{\beta}_1$ is the slope coefficient from that regression. We can extend this argument to show that Y can be residualized as well.

Remark. If we think about this as if X_1 is schooling and X_2 is family income, we are essentially including only the uncorrelated parts of those variables through this method, and don't need to worry about the correlation between the two.

When thinking about *Omitted Variable Bias*, we can similarly characterize the projection of X_2 on X_1 by

$$0 = \mathbb{E}_n(X_1(X_2 - \tilde{\gamma}X_1)) = \mathbb{E}_n \left(X_1 \left(\hat{\beta}_2 X_2 - \hat{\beta}_2 \tilde{\gamma} X_1 \right) \right)$$

We can then substitute the first order condition from above, and get

$$0 = \mathbb{E}_n \left(X_1 \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_2 X_2 \right) \right) = \mathbb{E}_n \left(X_1 \left(Y - \hat{\beta}_1 X_1 - \hat{\beta}_2 \tilde{\gamma} X_1 \right) \right)$$

So we have the first order condition from regression Y on X_1 only. We can conclude that the slope coefficient

from that regression is

$$\tilde{\beta} = \hat{\beta}_1 + \tilde{\gamma}\hat{\beta}_2$$

If there is a *causal interpretation* of the projection of Y onto (X_1, X_2) , then the difference term $\tilde{\gamma}\hat{\beta}_2$ is (the sample analog of) the *omitted variable bias* incurred by omitting X_2 .

Remark. We haven't even really introduced the idea of random variables and expectation yet. The word *bias* here is really loose, and doesn't make sense in the world of projection. However, in the real economic word, it makes a lot of sense in most contexts – in that case, you are saying that the omitted variables systematically bias the estimator in some direction.

Example. Frisch-Waugh(-Lovell) General Statement We can now do the same thing, more generally. Partition as follows:

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \quad ; \quad \mathbb{E}XX' \equiv Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} = \begin{pmatrix} \mathbb{E}X_1X_1' & \mathbb{E}X_1X_2' \\ \mathbb{E}X_2X_1' & \mathbb{E}X_2X_2' \end{pmatrix}$$

$$\mathbb{E}XY \equiv Q_{XY} = \begin{pmatrix} Q_{1Y} \\ Q_{2Y} \end{pmatrix} = \begin{pmatrix} \mathbb{E}X_1Y \\ \mathbb{E}X_2Y \end{pmatrix}$$

and use notation \hat{Q}_{11} (etc) for sample analogs. Then, it can be shown that

$$\hat{\beta} = \begin{pmatrix} \left(\hat{Q}_{11} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{21} \right)^{-1} \left(\hat{Q}_{1Y} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{2Y} \right) \\ \left(\hat{Q}_{22} - \hat{Q}_{21}\hat{Q}_{11}^{-1}\hat{Q}_{12} \right)^{-1} \left(\hat{Q}_{2Y} - \hat{Q}_{21}\hat{Q}_{11}^{-1}\hat{Q}_{1Y} \right) \end{pmatrix}$$

Regressing X_1 on X_2 would yield coefficients $\hat{\gamma} \equiv \hat{Q}_{22}^{-1}\hat{Q}_{21}$ and residuals $\hat{\eta} \equiv X_1 - X_2\hat{Q}_{22}^{-1}\hat{Q}_{21}$. Hence, we have that

$$\begin{aligned} \mathbb{E}_n\hat{\eta}^2 &= \mathbb{E}_nX_1^2 + \hat{Q}_{12}\hat{Q}_{22}^{-1}\mathbb{E}_nX_2^2\hat{Q}_{22}^{-1}\hat{Q}_{21} - 2\mathbb{E}_nX_1X_2\hat{Q}_{22}^{-1}\hat{Q}_{21} \\ &= \hat{Q}_{11} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{21} \end{aligned}$$

By similar algebra, $\mathbb{E}_n\hat{\eta}Y = \hat{Q}_{1Y} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{2Y}$. It follows that if we projected Y on the residual from regressing X_1 on X_2 , we would get coefficient $\hat{\beta}$, where

$$\tilde{\beta}_1 = \left(\hat{Q}_{11} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{21} \right)^{-1} \left(\hat{Q}_{1Y} - \hat{Q}_{12}\hat{Q}_{22}^{-1}\hat{Q}_{2Y} \right) = \hat{\beta}_1$$

Remark. Verbally, the multivariate OLS coefficient on X_1 is the coefficient one would get by regressing Y on the residual from regressing X_1 on all other covariates. We can show that this statement remains true if we replace Y with the residual $Y - \mathcal{P}_{X_2}(Y)$, where \mathcal{P}_{X_2} is the population projection onto X_2 . This may look like a curiosity now, but is an important starting point for partially linear models and other such things. An immediate payoff is that if you recall it, you'll never forget the own-variance formula for multivariate regression coefficients.

There are two ways to interpret OLS. Neither of them are uniquely 'right', but you should always be clear about which one you are appealing to. The *Best Linear Prediction* is an interpretation that makes sense under extremely general conditions. It's the notion of a linear model that is generalized in most predictive (notably, data science / statistical learning) applications. The *Causal (or Structural) Linear Model* is more demanding, but allows for causal interpretation. It's the notion of linear model that is generalized in most causal (for example, Instrumental Variables) applications. Note that this does not preclude predictive application of linear models as a component of causal inference. A salient example is the "first stage" in IV regression.

Best Linear Prediction. Write $Y = m(X) + \varepsilon$, where $m(x) \equiv \mathbb{E}(Y \mid X = x)$. That $\mathbb{E}(\varepsilon \mid X) = 0$ is now a tautology. We can show that $b^* \equiv (\mathbb{E}XX')^{-1}\mathbb{E}(XY)$, if it exists, minimizes $\mathbb{E}(Y - X'b)^2$. That

is, $\hat{Y} \equiv \mathcal{P}_X Y \equiv X'b^*$ is the *best linear predictor under square loss*. Furthermore, under these conditions, $\mathbb{E}Xe = 0$, where $e \equiv Y - \mathcal{P}_X(Y)$, meaning that the projection error e is not correlated with X .

Theorem 1.1. *If a weak law of large numbers applies to both $\frac{1}{n} \sum_{i=1}^n X_i X_i'$ and $\frac{1}{n} \sum_{i=1}^n X_i Y_i$ and $\mathbb{E}XX'$ is nonsingular, then b^* is uniquely defined and $\hat{\beta} \xrightarrow{P} b^*$.*

Remark. The best linear predictor \hat{Y} is uniquely defined even if $\mathbb{E}XX'$ is singular. In that case, b^* is not unique. This is really important for models with lots of covariates – think in statistical learning / machine learning.

The (Causal/Structural) Linear Model Write $Y = X'\beta + \varepsilon$, where $\mathbb{E}(\varepsilon | X) = 0$. Equivalently, $m(x) = \mathbb{E}(Y | X = x) = x'\beta$. In this version, $\mathbb{E}(\varepsilon | X) = 0$ is *not* tautological! Our assumptions become much stronger. To be precise, we are assuming that (i) ε is mean-independent of X , and (ii) the mean of ε is zero. These are new assumptions! We pay a large cost, but there are also large benefits. This model allows for *causal interpretation* of the estimand: in expectation, a change ΔX causes a corresponding change $\Delta X'\beta$ in Y . Importantly, this difference isn't just about interpretation – some important results are only available under the stronger assumptions.

Remark. Remember that within limits, a linear model can capture nonlinearities. Some examples:

1. Polynomial expansion:

$$Y = \beta_0 + \beta_1 H + \beta_2 H^2 + \dots + \varepsilon$$

2. Log or log-log regression:

$$\ln Y = \ln A + \alpha \ln K + (1 - \alpha) \ln L + \dots + \varepsilon$$

(Note! This changes the necessary assumption on ε , as in the primal it is now multiplied by the covariates)

3. Treatment effects with interactions:

$$Y = \beta + \delta \cdot \text{treatment} + \gamma \cdot \text{female} \cdot \text{treatment} + \dots + \varepsilon$$

Indeed, many “big data” models are high-dimensional but linear! Think of basically any machine learning context.

OLS as a Random Variable. The central questions are: What can we say about the estimator as a random variable? Are there conditions under which it has desirable properties, notably if our objective is to learn about β (or possibly the population $\mathcal{P}_X(Y)$)?

Consider drawing n samples and taking β_i , for $i \in \{1, \dots, n\}$ to be a random variable. If we show the histogram (in the slides) of these $\{\beta_i\}$, we can see that as n increases they tend to look normal! However, this is not necessarily a central limit theorem. Actually, the distribution of the estimators will approach the distribution of the error terms as n increases, under weaker assumptions. However, we assume normally distributed errors a lot.

Remark. We had a small aside to wonder whether we think of the vertical squared distance between each point and the line, or the shortest distance between the point and the line squared. The second is actually different from OLS – it's precisely *principal component analysis*!

Remark. What if we minimized the horizontal distance? Then we would get the projection of X onto Y rather than *vice versa* – called *reverse regression*. This is expanded on a lot in the first homework – basically, the coefficients are normalized to the variances of X and Y respectively. They are the same if and only if $\text{Var}(X) = \text{Var}(Y)$.

Remark. What if we minimize absolute values instead of squares? That would be *median regression* – at a population level, the median will solve this problem. By using other loss functions, we could tease out the

other quantiles, and by using all of them we would get *quantile regression*, where we treat the quantiles as being heterogeneously treated.

For finite sample theory, we will make the following assumptions:

Assumption 1.1. *In data matrix notation, we have that*

1. $Y = X\beta + \varepsilon$ ‘linearity’
2. $\mathbb{E}(\varepsilon \mid X) = 0$ ‘strong exogeneity’
3. $\text{rank}(X) = K$ a.s., where $X \in \mathbb{R}^{n \times K}$ ‘rank condition’ (equivalent: $X'X$ nonsingular a.s.)
4. $\mathbb{E}(\varepsilon\varepsilon' \mid X) = \sigma^2 I_n$ ‘spherical error’

The first two assumptions together imply a causal linear model. Assumptions that are natural when considering OLS as the ‘best linear predictor’ do not suffice to attain unbiasedness of $\hat{\beta}$. These further imply that ε is zero in expectation conditional on *all* covariates – including past and future realizations. That’s quite strong, but matters a lot more in time series econometrics. If we assume the data are i.i.d., this is not stronger than the earlier assumption that they are independent vector-wise.

The third assumption is an identification condition, and will fail if any covariates are linearly dependent on each other. We already talked about this previously, but note that in finite samples we are immediately excluding discrete covariates! Though the probability may be small, it doesn’t meet the criterion for almost surely. We may describe $\hat{\beta}$ as ‘conditionally unbiased’, where we are conditioning on X .

Assumption four combines conditional uncorrelatedness and homoskedasticity of errors. The latter makes sense only in the causal model because, even if the true regression error $\varepsilon = Y - m(X)$ is homoskedastic, the projection error

$$e = Y - \mathcal{P}_X(Y) = Y - m(X) + m(X) - \mathcal{P}_X(Y) = \varepsilon + m(X) - \mathcal{P}_X(Y)$$

is not.

We have a hidden assumption that $\mathbb{E}\|X\|^2 < \infty$. This is an existence result, and not so strong generally – if we assume X is non-stochastic, nothing here relies on it. It is an assumption, however.

These assumptions give us the following theorems:

Theorem 1.2. Finite Sample Bias and Variance *Under Assumptions 1.1, we have that*

1. $\mathbb{E}(\hat{\beta} \mid X) = \beta$
2. $\text{Var}(\hat{\beta} \mid X) = \sigma^2(X'X)^{-1}$

i.e. $\hat{\beta}$ is unbiased and its variance is determined.

Proof. We first observe that

$$\hat{\beta} = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X\beta + \varepsilon) = \beta + \underbrace{(X'X)^{-1}X'\varepsilon}_{= \text{estimation error}}$$

Therefore, the first two claims are that

$$\begin{aligned} \mathbb{E}\left((X'X)^{-1}X'\varepsilon \mid X\right) &= 0 \\ \text{Var}\left((X'X)^{-1}X'\varepsilon \mid X\right) &= \sigma^2(X'X)^{-1} \end{aligned}$$

We have directly that

$$\begin{aligned}\mathbb{E}\left((X'X)^{-1}X'\varepsilon \mid X\right) &= (X'X)^{-1}X'\underbrace{\mathbb{E}(\varepsilon \mid X)}_{=0} = 0 \\ \text{Var}\left((X'X)^{-1}X'\varepsilon \mid X\right) &= (X'X)^{-1}X'\sigma^2I_nX'(X'X)^{-1} = \sigma^2(X'X)^{-1}\end{aligned}$$

□

Theorem 1.3. Gauss-Markov Theorem Under Assumptions 1.1, if an estimator $\tilde{\beta}$ is linear (in Y) and unbiased, then

$$\text{Var}(\tilde{\beta} \mid X) \geq \sigma^2(X'X)^{-1}$$

Proof. We assume that $\tilde{\beta} = CY$ for some C (that may depend on X). Define $D \equiv C - (X'X)^{-1}X'$, and we have that

$$\begin{aligned}\beta &= \mathbb{E}\left(((X'X)^{-1}X' + D)Y \mid X\right) && \text{(Unbiased)} \\ &= \beta + \mathbb{E}(DY \mid X) \\ &= \beta + \mathbb{E}(D(X\beta + \varepsilon) \mid X) \\ &= \beta + \mathbb{E}(DX\beta \mid X) + D\underbrace{\mathbb{E}(\varepsilon \mid X)}_{=0} \\ &\implies \mathbb{E}(DX\beta \mid X) = 0\end{aligned}$$

This result is only possible if conditional on X the expression $DX\beta$ is non-stochastic – we have that $DX\beta = 0$ for any β ! This holds only if $DX = 0$. So finally, we have

$$\begin{aligned}\text{Var}(\tilde{\beta} \mid X) &= \text{Var}\left(((X'X)^{-1}X' + D)Y \mid X\right) \\ &= \text{Var}\left(((X'X)^{-1}X' + D)(X\beta_0 + \varepsilon) \mid X\right) \\ &= \text{Var}\left(((X'X)^{-1}X' + D)\varepsilon \mid X\right) \\ &= \sigma^2((X'X)^{-1}X' + D)((X'X)^{-1}X' + D)' \\ &= \sigma^2\left((X'X)^{-1}X'X(X'X)^{-1} + DX(X'X)^{-1} + (X'X)^{-1}X'D' + DD'\right) \\ &\geq \sigma^2(X'X)^{-1} = \text{Var}(\hat{\beta})\end{aligned}$$

where the conclusion follows from the fact that DD' is positive semi-definite. □

Question. Is homoskedasticity necessary for this result?

Answer. Yes! Consider the case where $\mathbb{E}\varepsilon\varepsilon' = \Omega$ is known and diagonal but its diagonal entries are not the same. Then the Gauss-Markov assumptions apply to the transformed model

$$\Omega^{-\frac{1}{2}}Y = \Omega^{-\frac{1}{2}}X + \Omega^{-\frac{1}{2}}\varepsilon$$

so the estimator

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

is the best linear unbiased estimator. However, $\hat{\beta}_{WLS} = \hat{\beta}$ if and only if $\Omega = \sigma^2I_n$ for some σ^2 . This is called the *Weighted Least Squares* estimator, which in this case would perform better than OLS.

Some important closed-form expressions Consider a simple linear regression $Y = \alpha + \beta X + \varepsilon$. Then we have that:

$$\hat{\beta} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})Y_i}{\sum_{i=1}^n (X_i - \bar{X})X_i}$$

$$\text{Var}(\hat{\beta} | X) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

where σ^2 is the variance of ε . In a multivariate regression, the variance of the k th component of $\hat{\beta}$ is

$$\text{Var}(\hat{\beta}_k | X) = \frac{\sigma^2}{(1 - R_k^2) \sum_{i=1}^n (X_{ki} - \bar{X}_k)^2}$$

where R_k^2 is the R^2 of the regression of X_k on the other covariates. Why is this obvious? Frisch-Waugh-Lovell, of course! The factor $1/(1 - R_k^2)$ is sometimes called the *variance inflation factor (VIF)*.

Remark. Know the expression for the variance of the k th component of $\hat{\beta}$ by heart! You should internalize it from Frisch-Waugh-Lovell.

We also have sample analogs and estimators for σ^2 : The *sample analog* (which is a *method of moments estimator*) of σ^2 is

$$\hat{\sigma}^2 \equiv \frac{1}{n} \hat{\varepsilon}' \hat{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

However, we can also show (exercise?) that $\mathbb{E}(\hat{\sigma}^2 | X) = \frac{n-K}{n} \sigma^2$. Why? Heuristically, the random variable $\hat{\varepsilon}$ has only $(n - K)$ degrees of freedom because it is constrained by the K equations $X' \hat{\varepsilon} = 0$. It is more common to use the unbiased

$$s^2 \equiv \frac{1}{n - K} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

Estimators of the standard deviation are often important. We will call them *standard errors*. $\sqrt{s^2}$ is the standard error of the regression, and $SE(\hat{\beta}) = (s^2 [(X'X)^{-1}]_{kk})^{1/2}$ is the standard error of $\hat{\beta}_k$.

Remark. We **cannot** claim that these estimators are unbiased! Further, the use of ‘standard error’ is dominant in econometrics, but other disciplines¹ ‘standard error’ and ‘[sampling] standard deviation’ may be synonyms. In this case, we use ‘*estimated standard errors*’.

Remark. The spherical error assumption was only (fully) used for the variance expressions (and Gauss-Markov). Recall that from the algebra:

$$\begin{aligned} \hat{\beta} &= \beta + \underbrace{(X'X)^{-1} X' \varepsilon}_{\text{estimation error}} \\ \implies \text{Var}(\hat{\beta} | X) &= (X'X)^{-1} X' \underbrace{\mathbb{E}(\varepsilon \varepsilon' | X)}_{=D} X (X'X)^{-1} \end{aligned}$$

Spherical error ($D = \sigma^2 I_n$) leads to simplification, but as long as we can estimate D , it is not necessary. Recall that

$$\text{Var}(\hat{\beta} | X) = (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \varepsilon_i^2 \right) (X'X)^{-1}$$

Assuming only heteroskedasticity (*i.e.* D remains diagonal), we have the *oracle estimator*

$$\hat{\text{Var}}_{\text{oracle}}(\hat{\beta} | X) \equiv (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \varepsilon_i^2 \right) (X'X)^{-1}$$

¹Statistics...

This is an unbiased estimator, but is not available. Plugging in $\hat{\varepsilon}_i$ leads to a plausible estimator:

$$\hat{\text{Var}}_{\text{HC0}}(\hat{\beta} \mid X) \equiv (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{\varepsilon}_i^2 \right) (X'X)^{-1}$$

This is biased, which motivates a degree of freedom adjustment:

$$\hat{\text{Var}}_{\text{HC1}}(\hat{\beta} \mid X) \equiv \frac{n}{n-K} (X'X)^{-1} \left(\sum_{i=1}^n X_i X_i' \hat{\varepsilon}_i^2 \right) (X'X)^{-1}$$

There are a ton of similar estimators. HC0 is the original (Eicker-White) heteroskedasticity robust variance estimator, and HC1 is the industry standard (*i.e.* it is the STATA default). Neither is obviously best, and Hansen has some other options. In applied work, these are dominant because homoskedasticity is rarely plausible. Using weighted least squares is technically possible, but rarely used because it requires knowing *ex ante* the structure of the heteroskedasticity. Clustered standard errors are similar to this, though we omit them for now.

From here on, we will make the following (strong) assumption:

Assumption 1.2. ε has a normal distribution:

$$(\varepsilon \mid X) \sim \mathcal{N}(0, \sigma^2 I_n)$$

With this, we have the following:

Theorem 1.4. Define $s^2 = \frac{(Y-X\beta)'(Y-X\beta)}{n-K}$ and let the matrix $R \in \mathbb{R}^{r \times K}$ have maximal rank $r \leq K$. Under Assumptions 1.1 and 1.2, we then have

$$(\hat{\beta} - \beta) \mid X \sim \mathcal{N}(0, \sigma^2 (X'X)^{-1})$$

and:

$$\begin{aligned} t\text{-ratio} = t &\equiv \frac{\hat{\beta}_k - \beta_k}{\left(s^2 [(X'X)^{-1}]_{kk} \right)^{1/2}} \sim t_{n-K} \\ F\text{-statistic} = F &\equiv \frac{(R\hat{\beta} - R\beta)'(R(X'X)^{-1}R')^{-1}(R\hat{\beta} - R\beta)}{s^2 r} \sim F_{r, n-K} \end{aligned}$$

Thus, the null hypothesis $\mathbb{H}_0 : R\beta = r$ can be tested with exact size control by comparing

$$\frac{(R\hat{\beta} - r)'(R(X'X)^{-1}R')^{-1}(R\hat{\beta} - r)}{s^2 r}$$

to the relevant quantiles of $F_{r, n-K}$, etc.

Proof. The first part, from section: We showed earlier that $\hat{\beta} = \beta + (X'X)^{-1}X'\varepsilon$, so

$$\mathbb{E}[\hat{\beta} - \beta \mid X] = \mathbb{E}[\hat{\beta} \mid X] - \beta = \beta + (X'X)^{-1}X' \underbrace{\mathbb{E}[\varepsilon \mid X]}_{=0} - \beta = 0$$

and

$$\text{Var}(\hat{\beta} - \beta \mid X) = \text{Var}(\hat{\beta} \mid X) = \sigma^2 (X'X)^{-1}$$

Conclusion follows by also noting that the sum of normally-distributed random variables.

The rest follows directly from the definitions of the relevant distributions, and our assumptions on R . \square

Remark. Recall that the t -distribution converges quickly to the standard normal as degrees of freedom

increase, but the quantiles (and thus associated critical values) can be quite different under small degrees of freedom.

2 The Linear Model in Large Samples

We will now look at large sample (asymptotic) theory; results will be weaker in that they maximally only hold *almost surely* (i.e. we care about stochastic convergence). However, our assumptions are also a lot weaker:

Assumption 2.1. (X, Y) are i.i.d., $\mathbb{E}Y^2 < \infty$ and $\mathbb{E}\|X\|^2 < \infty$, and $\mathbb{E}(XX')$ is positive definite.

Remark. These assumptions suffice for the projection coefficient to be well-defined: $b^* \equiv (\mathbb{E}XX')^{-1}\mathbb{E}XY$. We have:

Theorem 2.1. Under Assumptions 2.1, we have:

$$\hat{\beta} \equiv (\mathbb{E}_n XX')^{-1} \mathbb{E}_n XY \xrightarrow{P} b^*$$

Proof. By the weak law of large numbers, we have that $\mathbb{E}_n XX' \xrightarrow{P} \mathbb{E}XX'$ and $\mathbb{E}_n XY \xrightarrow{P} \mathbb{E}XY$. By the Continuous Mapping Theorem, it follows that $(\mathbb{E}_n XX')^{-1} \xrightarrow{P} (\mathbb{E}XX')^{-1}$, since as n increases $\mathbb{E}_n XX'$ is nonsingular with probability approaching 1. The claim follows directly from Slutsky's Theorem. \square

We have not yet used any assumptions that set the structural linear model apart from the best linear predictor under square loss interpretation, but we've made no causal claims. However, the asymptotics apply to both!

We must assume the following to make causal claims, if the linear model $Y = X'\beta + \varepsilon$ is maintained:

Assumption 2.2. $\mathbb{E}X\varepsilon = 0$, also called [predetermination](#). Additionally, we strengthen the moment assumption: $\mathbb{E}Y^4 < \infty$ and $\mathbb{E}\|X\|^4 < \infty$.

Remark. This is significantly weaker than the unbiasedness assumption – we only need that ε is uncorrelated with the contemporaneous regressors. This ensures that $b^* = \beta$ (and, by implication, consistency):

$$b^* = (\mathbb{E}XX')^{-1}\mathbb{E}XY = (\mathbb{E}XX')^{-1}\mathbb{E}X(X'\beta + \varepsilon) = \beta + (\mathbb{E}XX')^{-1}\mathbb{E}X\varepsilon = \beta$$

We could alternatively show consistency from scratch using this assumption, the way Hayashi does.

To look at the asymptotic distribution, note that

$$\begin{aligned} \hat{\beta} &= (\mathbb{E}_n XX')^{-1} \mathbb{E}_n XY = (\mathbb{E}_n XX')^{-1} \mathbb{E}_n X(Xb^* + e) \\ \implies \hat{\beta} - b^* &= (\mathbb{E}_n XX')^{-1} \mathbb{E}_n Xe \end{aligned}$$

where b^* is the population projection coefficient and e is the (true underlying) projection error. By predetermination, $(\mathbb{E}_n XX')^{-1} \mathbb{E}_n Xe \xrightarrow{P} 0$. However, there is the suggestion of an asymptotic result – it seems natural that

$$\sqrt{n}(\mathbb{E}_n Xe) \xrightarrow{d} \mathcal{N}(0, \Omega) = \mathcal{N}(0, \mathbb{E}(XX'e^2))$$

where the last equality just defines Ω . If that were the case, we would easily have that

$$\sqrt{n}(\hat{\beta} - b^*) \xrightarrow{d} \mathcal{N}(0, (\mathbb{E}XX')^{-1}\Omega(\mathbb{E}XX')^{-1})$$

This derivation is basically true.² We just need to ensure that all terms exist. For this, our stronger Assumptions 2.2 suffice. If the kurtosis exists, we can argue that Ω is finite, by repeatedly using Cauchy-Schwartz. For any one element of Ω ,

$$\begin{aligned} |\mathbb{E}(X_k X_\ell e^2)| &\leq \mathbb{E}|X_k X_\ell e^2| = \mathbb{E}(|X_k||X_\ell||e^2|) \\ &\leq (\mathbb{E}X_k^2 X_\ell^2)^{1/2} (\mathbb{E}e^4)^{1/2} \leq (\mathbb{E}X_k^4)^{1/4} (\mathbb{E}X_\ell^4)^{1/4} (\mathbb{E}e^4)^{1/2} < \infty \end{aligned}$$

²'Morally true' - Jörg

we finish by writing

$$\sqrt{n}(\hat{\beta} - b^*) \xrightarrow{d} (\mathbb{E}XX')^{-1}\mathcal{N}(0, \Omega) = \mathcal{N}(0, (\mathbb{E}XX')^{-1}\Omega(\mathbb{E}XX')^{-1})$$

where we again use Slutsky and the properties of normal distributions. Formally, we have:

Theorem 2.2. *With Assumptions 2.1 and Assumptions 2.2, we have that*

$$\begin{aligned} \sqrt{n}(\hat{\beta} - b^*) &\xrightarrow{d} \mathcal{N}(0, \text{aVar}(\hat{\beta})) \\ \text{where } \text{aVar}(\hat{\beta}) &= (\mathbb{E}XX')^{-1}\Omega(\mathbb{E}XX')^{-1} \\ &[= Q_{XX}^{-1}\Omega Q_{XX}^{-1} = \Sigma_{XX}^{-1}\Omega\Sigma_{XX}^{-1}] \end{aligned}$$

Proof. A generalization of above. □

Remark. Note that we are *not* assuming a linear model here! We actually get this result under relatively limited assumptions, we only need that the moment conditions exist. If we want to get inference results, we need more, but as a projection result this still holds.

Definition. aVar is the *asymptotic variance*, the variance of the limiting distribution. In general, this is not necessarily the asymptotic limit of an estimator's squared variance, though the current assumptions suffice.

Remark. This theorem provides the *joint* asymptotic distribution of estimates. The information contained in joint normality is relevant for:

1. Inference on a linear combination of estimates, *e.g.* their sum or difference. This could also be achieved by reparameterization, but that's impractical.
2. Joint inference, *i.e.* confidence ellipsoids, on several coefficients
3. Inference on a known, differentiable function of β through the Delta method (conceptually straightforward, but very important in practice! See the textbook for an example worked through)
4. Conservative inference on a known, nondifferentiable function of β through projection (*i.e.* operate the function on ever b in the confidence ellipsoid). In structured cases, you may be able to improve on this – ask Jörg if this question arises in your research!

Remark. Results in this course hold pointwise as $n \rightarrow \infty$ for *given* parameter values, not *uniformly* over parameter values. How big of a problem is this? With some more effort, most results in this course are available uniformly in ‘nice’ cases. However, there are several cases that are not nice and are empirically relevant: (i) estimators that can be corner solutions of their problems, (ii) estimation of maxima, (iii) rare events, and (iv) post-model selection estimation and inference. In these, pointwise perspectives can be quite misleading.

Theorem 2.3. *Under Assumptions 2.1 and Assumptions 2.2, we have that*

$$\text{a}\hat{\text{Var}}_{HC0} \equiv (\mathbb{E}_n XX')^{-1} \hat{\Omega} (\mathbb{E}_n XX')^{-1} \xrightarrow{P} \text{aVar}(\hat{\beta})$$

where $\hat{\Omega} := \mathbb{E}_n[XX'\hat{e}^2]$, and similarly for $HC1$, and so on. Basically, all of these are consistent.

Proof. (Sketch) The bottleneck is the consistency of $\hat{\Omega}$. Write:

$$\hat{\Omega} = \frac{1}{n} \sum_{i=1}^n X_i X_i \hat{e}_i^2 = \underbrace{\frac{1}{n} \sum_{i=1}^n X_i X_i' e_i^2}_{\xrightarrow{P} \Omega} + \underbrace{\frac{1}{n} \sum_{i=1}^n X_i X_i (\hat{e}_i^2 - e_i^2)}_{\xrightarrow{P} 0}$$

That the rightmost term disappears can be shown through (tedious) repeated use of Cauchy-Schwartz and Hölder. □

Proof. (Assuming Homoskedasticity) If we assume that $\mathbb{E}(e^2 | X) = \sigma^2$, then we have the simplification

$$\text{aVar}(\hat{\beta}) = (\mathbb{E}XX')^{-1}\sigma^2$$

and showing the consistency of

$$\text{a}\hat{\text{Var}}(\hat{\beta}) = (\mathbb{E}XX')^{-1}s^2$$

is simple, and requires weaker assumptions (specifically, second moments suffices). However, recall that this makes sense only for a structural linear model (and is still restrictive). \square

Inference. Let $r : \mathbb{R}^k \rightarrow \mathbb{R}$ be a continuously differentiable function with $\nabla r(\cdot) = R(\cdot)$ (the easiest example is where r extracts a component of β). Define $\theta = r(\beta)$ and $\hat{\theta} = r(\hat{\beta})$. By Delta Method, standard convergence results, and Slutsky, we have

$$\begin{aligned} \sqrt{n}(\hat{\theta} - \theta) &\xrightarrow{d} \mathcal{N}(0, \text{aVar}(\hat{\theta})) \\ \text{aVar}(\hat{\theta}) &= R(\beta) \text{aVar}(\hat{\beta}) R(\beta)' \\ \text{a}\hat{\text{Var}}(\hat{\theta}) &\equiv R(\hat{\beta}) \text{aVar}(\hat{\beta}) R(\hat{\beta})' \xrightarrow{p} R(\beta) \text{aVar}(\hat{\beta}) R(\beta)' \\ \implies t(\theta) &\equiv \frac{\hat{\theta} - \theta}{SE(\hat{\theta})} \equiv \frac{\sqrt{n}(\hat{\theta} - \theta)}{\left(R(\hat{\beta}) \text{a}\hat{\text{Var}}(\hat{\beta}) R(\hat{\beta})'\right)^{1/2}} \xrightarrow{d} \mathcal{N}(0, 1) \end{aligned}$$

This is the asymptotic t -statistic. For this to hold, we need that $\text{aVar}(\hat{\theta})$ is finite. A sufficient condition is that $R(\beta) \neq 0$ and that $\text{aVar}(\hat{\beta})$ has full rank.

Definition. Dividing by the standard error like this is called *studentization*. It ensures that the asymptotic distribution does not depend on unknown parameters, ensuring that the t -statistic (and others) are *asymptotic pivots*.

The previous result lets us create hypothesis tests and confidence intervals where the asymptotic sizes converge. Let $\Phi(\cdot)$ denote the standard normal cdf and define the quantiles $\Phi^{-1}(1 - \alpha) := c_\alpha$. Then

$$\begin{aligned} \mathbb{P}\{|t(\theta)| \leq c_{\alpha/2}\} &\xrightarrow{p} 1 - \alpha \\ \mathbb{P}\{t(\theta) \in CI_\alpha(\theta)\} &\xrightarrow{p} 1 - \alpha \\ CI_\alpha(\theta) &= \left[\hat{\theta} - c_{\alpha/2} \cdot SE(\hat{\theta}), \hat{\theta} + c_{\alpha/2} \cdot SE(\hat{\theta})\right] \end{aligned}$$

and we compute one-sided confidence intervals similarly. If $r(\beta) = p'\beta$ for some known vector p , then $R(\cdot) = p'$ and the t -statistic simplifies to

$$t(\theta) \equiv \frac{\hat{\theta} - \theta}{SE(\hat{\theta})} = \frac{\sqrt{n}(\hat{\theta} - \theta)}{\left(p' \text{a}\hat{\text{Var}}(\hat{\beta}) p\right)^{1/2}}$$

If p is a basis vector, this further simplifies the t -statistic for an individual coefficient. Choices like $p = (0, 1, -1, 0, \dots, 0)$ lets us test the equality of two coefficients.

Under the causal linear model, another application is to $p = x$, in which case $\theta = \mathbb{E}(Y | X = x)$, which is called a *regression interval*. Note that the standard error depends on x and will be smaller for more central values of x .

Remark. The regression interval is *not* a forecast confidence interval! For forecast intervals, we must take ε_t into account.

We can generalize this to $\theta = r(\beta)$ where $r : \mathbb{R}^k \rightarrow \mathbb{R}^q$ is vector-valued. With this, $R(\cdot)$ is now the Jacobian

of r , we have that

$$\begin{aligned}
\sqrt{n}(\hat{\theta} - \theta) &\xrightarrow{d} \mathcal{N}(0, \text{aVar}(\hat{\theta})) \\
\text{aVar}(\hat{\theta}) &= R(\beta) \text{aVar}(\hat{\beta}) R(\beta)' \\
\text{a}\hat{\text{Var}}(\hat{\theta}) &\equiv R(\hat{\beta}) \text{aVar}(\hat{\beta}) R(\hat{\beta})' \xrightarrow{p} R(\beta) \text{aVar}(\hat{\beta}) R(\beta) \\
\Rightarrow W(\theta) &\equiv \sqrt{n}(\hat{\theta} - \theta)' (\text{a}\hat{\text{Var}}(\hat{\theta}))^{-1} \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \chi_q^2
\end{aligned}$$

This is the *Wald* statistic. It can induce confidence regions similarly to the confidence intervals above – however these are ellipsoids. These will be appropriate if one is genuinely interested in simultaneous inference of several scalars, and their projections onto the axes will be valid but (possibly very) conservative confidence intervals.

Remark. A sufficient condition for this to hold is that both $R(\beta)$ and $\text{aVar}(\hat{\beta})$ are full-rank, which would mean that the hypotheses must be (locally) linearly independent at the truth – for linear hypotheses, this is easy to check and a global property.

3 Instrumental Variables

Remark. Instrumental variables are of huge importance across economics, and are one of the greatest contributions of econometrics to empirical methods across science. They are actively used in causal inference across disciplines, specifically in biostatistics. Their appeal is that they allow for causally interpretable estimates if we think that (i) the linear model (or generalizations) is structural so β has causal interpretation, but (ii) ε correlates with X , because it absorbs relevant omitted variables. Of course, this remarkable result requires strong assumptions.

Model. For simplicity, consider simple linear regression

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where we cannot assume that $\mathbb{E}X\varepsilon = 0$. The interpretation is that β_1 has a causal interpretation, but we could not estimate all relevant covariates. Now suppose that we have a random variable Z with the following properties:

$$\underbrace{\text{cov}(Z, X) \neq 0}_{\text{Relevance}} \quad \text{and} \quad \underbrace{\text{cov}(Z, \varepsilon) = 0}_{\text{Validity}}$$

This implies that

$$\frac{\text{cov}(Z, Y)}{\text{cov}(Z, X)} = \frac{\beta_1 \text{cov}(Z, X) + \text{cov}(Z, \varepsilon)}{\text{cov}(Z, X)} = \beta_1$$

Example. (Ginburgh & van Ours (2003)) Y is the career success of a classical musician, X is their placement in a prestigious competition for young musicians, and Z is the order of appearance at the competition. The effect of X on Y is interesting, but ‘talent’ impacts both. However, appearing late in the competition (verifiably) predicts success! Since order of appearance is random, it serves as an *instrumental variable* (or just *instrument*).

Definition. *Relevance* requires that $\text{cov}(Z, X) \neq 0$. We need there to be some instrument-induced variation to play around with. Otherwise, we could pay a research assistant to flip coins all day and use that as an instrument.

Relevance must be **testable**. The covariance is consistently estimated with its sample analog. Indeed, it is standard practice to report the F -statistic from a ‘first-stage regression’ of X on Z .

Definition. *Validity* requires that $\text{cov}(Z, \varepsilon) = 0$. We need the instrument-induced variation to be exogenous. Otherwise, we could just use X as an instrument for itself.

Validity is **not testable**. This will change once we have more instruments than regressors.

Definition. We can illustrate causal models using *Directed Acyclic Graphs (DAG)*. Arrows read as ‘causes...’. The traditional OLS graph is Figure 2.

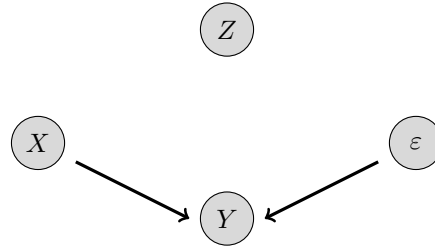


Figure 2: OLS DAG. The r.v. Z plays no role, and $\text{cov}(X, \varepsilon) = 0$.

If X is *endogenous*, the DAG may look like Figure 3.

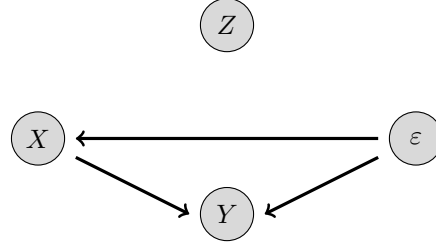


Figure 3: OLS DAG. The r.v. Z plays no role, but $\text{cov}(X, \varepsilon) \neq 0$.

The typical DAG representation of instrumental variables is Figure 4.

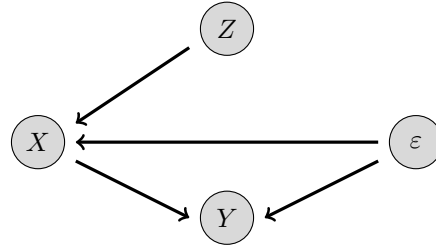


Figure 4: IV DAG. $\text{cov}(Z, X) \neq 0$, $\text{cov}(Z, \varepsilon) = 0$, and $\text{cov}(X, \varepsilon) \neq 0$.

Example. Instrumental variables in medical statistics: *Mendelian randomization* is an instrumental variables technique. The idea is that, in our terminology, genetic variation exogenous. Think about Y = coronary disease, X = HDL cholesterol level, and ε = confounders (of which there are many – diet and lifestyle heavily affect HDL cholesterol level). If we have some genetic variant Z that increases HDL cholesterol level, we can use that as an instrument. The bottleneck assumption is that Z affects *only* HDL cholesterol level, not coronary disease.

Example. Some famous³ examples:

1. [Wright \(1928\)](#) was the first example – he wanted to estimate supply and demand elasticities for vegetable oils, but realized that quantities and prices are equilibrium quantities and prices, so endogenous. He used demand shifters (price changes of substitutes) and supply shifters (weather) as instruments. He averaged the estimators, which is something we would do now.
2. [Angrist & Krueger \(1991\)](#) aim to estimate the returns to compulsory schooling, but selection into schooling is correlated. They instrument it by quarter of birth, where a different quarter may lead to another year ‘exposed’ to compulsory schooling.
3. [Angrist & Evans \(1998\)](#) want to estimate the effect of having children on female labor force participation, but the decision to have children is endogenous. They instrument with gender of children – this is clearly random, but some parents prefer to have children of both genders so if they have two of the same, they are more likely to have a third than if the two are different genders.
4. [Card \(1993,5\)](#) looks to estimate the returns to college education, but college attendance is endogenous. He instruments college attendance with distance from a college growing up, which may be pivotal in attendance decisions.

Definition. We will derive the IV estimator as a *Generalized Method of Moments (GMM)* estimator. Recall

³Note! Famous \neq good!

the method of moments: If we assume

$$Y = X'\beta + \varepsilon, \mathbb{E}X\varepsilon = 0 \implies \mathbb{E}(X(Y - X'\beta)) = 0$$

then a natural idea for estimating β is to solve the empirical moment condition

$$\mathbb{E}_n(X(Y - X'\hat{\beta})) = 0$$

which clearly yields the OLS estimator (as the above is the FOC of the minimization problem). However, if we assume that $\mathbb{E}X\varepsilon \neq 0$, then OLS will be inconsistent for β because it will estimate the projection coefficient

$$b^* = (\mathbb{E}XX')^{-1}\mathbb{E}XY = \beta + (\mathbb{E}XX')^{-1}\mathbb{E}X\varepsilon$$

However, if we also observe a random k -vector Z with $\mathbb{E}Z\varepsilon = 0$ and $\text{rank}(\mathbb{E}ZX') = k$, then we have the moment condition

$$\mathbb{E}(Z(Y - X'\beta)) = 0 \implies \beta = (\mathbb{E}ZX')^{-1}\mathbb{E}ZY$$

where the assumptions assume that this is well-defined. Naturally, the estimator is

$$\hat{\beta}_{IV} = (\mathbb{E}_nZX')^{-1}\mathbb{E}_nZY$$

Note that

$$\hat{\beta}_{IV} = (\mathbb{E}_nZX')^{-1}\mathbb{E}_nZY = (\mathbb{E}_nZX')^{-1}\mathbb{E}_nZ(X'\beta + \varepsilon) = \beta + \underbrace{(\mathbb{E}_nZX')^{-1}\mathbb{E}_nZ\varepsilon}_{= \text{estimation error}}$$

From here, consistency follows essentially directly as $\mathbb{E}_nZ\varepsilon \xrightarrow{p} 0$. However, we cannot claim unbiasedness, even if we assume the stronger that $\mathbb{E}(\varepsilon | Z) = 0$. We will defer asymptotic theory until later. Note that some components of Z may also appear in X , because if some elements of X are unassociated with ε , they can act as their own instruments. By setting $Z = X$, we consider OLS as the special case of X instrumenting itself.

In the simple linear model, when X and Z are scalars, the IV slope estimator can be expressed as

$$\frac{\sum_{i=1}^n (Z_i - \bar{Z})(Y_i - \bar{Y})}{\sum_{i=1}^n (Z_i - \bar{Z})(X_i - \bar{X})}$$

The data matrix expression is, naturally, $\hat{\beta}_{IV} = (Z'X)^{-1}Z'Y$. For an application, consider regressing X on Z and then Y on \hat{X} . In data matrix notation, we have

$$\begin{aligned} \tilde{\beta} &= (\hat{X}'\hat{X})^{-1}\hat{X}'Y \\ &= ((Z'(Z'Z)^{-1}Z'X)'(Z'(Z'Z)^{-1}Z'X))^{-1}(Z'(Z'Z)^{-1}Z'X)'Y \\ &= (X'Z(Z'Z)^{-1}Z'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \\ &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y \\ &= (Z'X)^{-1}Z'Z(X'Z)^{-1}X'Z(Z'Z)^{-1}Z'Y \\ &= (Z'X)^{-1}Z'Y = \hat{\beta}_{IV} \end{aligned}$$

Remark. The IV estimator can be thought of as a two-stage model where we project X onto Z and then project Y onto the fitted values \hat{X} . This gives some straightforward intuition – we exploit only the variation in X that is due to⁴ variation in Z . This interpretation is closely related to the DAG interpretation, and it's why the regression of X on Z is often called the *first-stage regression*. It is usually reported and should be highly significant – a rule of thumb is that $F \geq 10$ for the overall first-stage regression.

⁴In a correlative sense, we have not made causal claims. We will return to this when thinking about machine learning and optimal instruments.

Also note that this does not require Z and X to be the same length – we can use more instruments than regressors. However, a caveat:

Definition. If $Z'X$ is not square (but assuming it still has maximal rank!) the algebraic deviation becomes

$$\begin{aligned}\tilde{\beta} &= (\hat{X}'\hat{X})^{-1}\hat{X}'Y \\ &= ((Z(Z'Z)^{-1}Z'X)'(Z(Z'Z)^{-1}Z'X))^{-1}(Z(Z'Z)^{-1}Z'X)'Y \\ &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y = \hat{\beta}_{TSLS}\end{aligned}$$

where we can go no further than the last step. This is the *two-stage least squares (TSLS)* estimator.

Remark. A major difference in the estimators is that for IV, $Z'\hat{\varepsilon} = Z'(Y - X\hat{\beta}) = 0$ by construction. Can this also be true in TSLS? No! With $\ell > k$ instruments, this is ℓ linear equations in k unknowns. This has some major implications: the estimator cannot set $Z'\varepsilon = 0$, but we can show that the estimator minimizes (in b) $\|Z'(Y - Xb)\|$ for some norm, not necessarily the best norm. However, the validity of instruments becomes testable – in large samples, we should have that $\frac{1}{n}Z'\hat{\varepsilon} \approx 0$. These considerations lead us from two-stage least squares to the generalized method of moments. Detailed asymptotic characterization of TSLS follows from GMM.

We can think about asymptotics of the simple instrumental variable case, where

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

where (Y, X, z) are scalars, and we make the further assumption of homoskedasticity.⁵ From previous algebra, we get that

$$\begin{aligned}\sqrt{n}(\hat{\beta}_1^{IV} - \beta_1) &= \frac{\sqrt{n}\frac{1}{n}\sum_{i=1}^n(Z_i - \bar{Z})\varepsilon_i}{\frac{1}{n}\sum_{i=1}^n(Z_i - \bar{Z})(X_i - \bar{X})} \\ &= \frac{\sqrt{n}\frac{1}{n}\sum_{i=1}^n(Z_i - \mathbb{E}Z)\varepsilon_i}{\frac{1}{n}\sum_{i=1}^n(Z_i - \mathbb{E}Z)(X_i - \mathbb{E}X)} + o_P(1) \\ &\xrightarrow{d} \frac{\mathcal{N}(0, \sigma_Z^2 \sigma^2)}{\rho \sigma_Z \sigma_X} = \mathcal{N}\left(0, \frac{\sigma^2}{\rho^2 \sigma_X^2}\right)\end{aligned}$$

We can compare this directly to

$$\sqrt{n}(\hat{\beta}_1^{OLS} - \beta_1) \xrightarrow{d} \mathcal{N}\left(0, \frac{\sigma^2}{\sigma_X^2}\right)$$

So the asymptotic variance of the IV estimator has the first stage explained variation in the denominator.

Remark. This is not a result about finite sample variance! In fact, the finite sample variance of the IV estimator does not exist.

A corollary of this result is that the IV estimator has higher asymptotic variance, as a trade-off to bias. If $\rho^2 = 1$, the two asymptotic variances coincide (in fact, the estimators algebraically coincide – this is equivalent to using X as an instrument for itself). The asymptotic variance diverges to ∞ as $\rho \rightarrow 0$, which suggests some intuition about weak instruments.

Definition. A *weak instrument* is an instrument Z where $\text{corr}(Z, X) \approx 0$, meaning that the instruments explain very little of the variation in X .

Example. (Following [Staiger & Stock \(1997\)](#)) To formally model a weak instrument, set $\rho_n = \frac{\rho}{\sqrt{n}}$. Asymptotic approximation is powerful – it allows us to invoke CLT and other results. However, in a pointwise perspective as it trivializes the problem as previous asymptotics hold for any $\rho \neq 0$. *Parameter drift* allows us to invoke asymptotic approximations without approximating away the problem. Compare to [Pitman Drift](#) for analyzing the local power of hypothesis tests. The idea is *not* that parameters actually change with n .

⁵We make this assumption because (i) the asymptotic variance is instructive, and (ii) it allows us to formally characterize weak instruments.

Rather, the idea is to internalize the intuition that whether an instrument is weak depends on ρ in relation to n . See [Goldberger](#) on [micronumerosity](#) for more on this.

We will develop this for scalar (X, Z) . See Hayashi for a general statement. The first- and second stage regressions are

$$\begin{aligned} Y &= \beta_0 + \beta_1 X + \varepsilon \\ X &= \gamma_0 + \frac{\gamma_1}{\sqrt{n}} Z + \eta \end{aligned}$$

Then

$$\hat{\beta}_1 - \beta_1 = \frac{\sqrt{n} \mathbb{E}_n(Z - \bar{Z}) \varepsilon}{\sqrt{n} \mathbb{E}_n(Z - \bar{Z})(X - \bar{X})}$$

but

$$\begin{aligned} \sqrt{n} \mathbb{E}_n(Z - \bar{Z})(X - \bar{X}) &= \sqrt{n} \mathbb{E}_n(Z - \bar{Z}) \left(\gamma_0 + \frac{\gamma_1}{\sqrt{n}} Z + \eta \right) \\ &= \underbrace{\gamma_1 \mathbb{E}_n(Z - \bar{Z}) Z}_{\xrightarrow{p} \gamma_1 \sigma_Z^2} + \sqrt{n} \mathbb{E}_n(Z - \bar{Z}) \eta \end{aligned}$$

Assuming CLT applies, we attain that

$$\hat{\beta}_1 - \beta_1 \xrightarrow{d} \frac{a}{\gamma_1 \sigma_Z^2 + b} \quad \text{where} \quad \begin{pmatrix} a \\ b \end{pmatrix} \sim \mathcal{N}(0, \Omega)$$

where Ω is the variance-covariance matrix of $(Z\varepsilon, Z\eta)$. This estimator is inconsistent! Indeed, it converges to a distribution. Observe that the extent of this problem scales inversely with $|\gamma_1|$ – as $|\gamma_1| \rightarrow \infty$, the problem vanishes.

Remark. One of the examples referenced above, [Angrist & Krueger \(1991\)](#), has an instrument that is arguably weak, though n is large. [Bound, Jaeger, & Baker \(1993\)](#) replicated some of their tables with a random ‘instrument’ they added to the data. This spawned a large literature in weak instruments. Here at Cornell, [Pepe](#) has done a lot of work on weak instruments.

Remark. We might think about using many instruments instead of one strong instrument. In fact, if we let the number of instruments grow as n grows, Hansen has a formalization of the fact that that estimator is inconsistent.

Remark. We can also encounter issues with ‘too-strong’ instruments. These do not exist in theory, but consider the example where (i) we think that X is endogenous, and (ii) $\text{corr}(Z, X) = 0.99$. Theoretically there’s no issue with this, as long as $\text{cov}(Z, \varepsilon) = 0$. However, that’s not testable, and intuitively it would be very strange if Z and X were so correlated and Z would not at all be covariant with the errors. So in practice, we often want $F \leq 25$, even though in theory higher F is better.

4 Generalized Method of Moments

The Generalized Method of Moments (GMM) and its relatives like the Method of Simulated Moments and Indirect Inference are of great importance in applied work. This statement (and the name!) are due to [Hansen \(1982\)](#), and contributed to his Nobel Prize. There were precursors to a large part of this theory, but we will develop a fairly general statement, though we restrict attention to linear moment conditions. Nonlinear GMM will be developed as a special case of extremum estimation.

We can think of GMM as extending TSLS in several ways:

1. Since we cannot set sample moments exactly to zero, we must choose a norm to minimize. Is there a best norm?
2. We allow for heteroskedasticity also in the estimation stage (Heteroskedasticity robust standard errors for TSLS are straightforward and standardly used, but we will see that in the estimation stage, TSLS can be argued to presume homoskedasticity.)
3. We consider testing instrument validity
4. It will become clear that restricting to linear moment conditions simplifies the math but is not essential

Remark. The *generalized* in GMM refers to the fact that we allow for (and explore the implications of!) overidentification, when we have more moment equations than parameters.

Definition. We know that

$$\mathbb{E}g(Y, X, Z; \theta) = 0$$

where $\theta \in \mathbb{R}^k$ and $g(\cdot)$ is a known smooth function mapping into \mathbb{R}^ℓ , with $\ell \geq k$. The case of $\ell > k$ will be called *overidentified*. We will assume $g(\cdot)$ is linear. This is not essential.

Remark. This yields OLS, IV, TSLS, and (after generalizing to multiple outcomes) seemingly unrelated regression equations (SURE) and simple panel data estimators as special cases. For future reference, consider also *probabilistic regression (probit)* where $\mathbb{E}[X(Y - \Phi(X'\beta))] = 0$, and best-response conditions such as Euler equations (this was the original application of GMM, from [Hansen & Singleton \(1983\)](#)).

The *GMM estimator* is

$$\begin{aligned}\hat{\theta}(W) &= \underset{\theta}{\operatorname{argmin}} J_n(\theta) \\ J_n(\theta) &= n\bar{g}_n(\theta)'W\bar{g}_n(\theta) \\ \bar{g}_n(\theta) &\equiv \frac{1}{n} \sum_{i=1}^n g(\theta; \cdot)\end{aligned}$$

where W is a weight matrix defining the norm that we minimize. Recall that if we are overidentified, we cannot set $J_n(\theta)$ to zero. We therefore have a family of estimators, and will think about how we optimally choose W . The scale factor in $J_n(\cdot)$ is for convenience, ensuring that it converges to a non-degenerate limit.

Example. We begin with the linear case, where

$$g(\beta, \cdot) = Z(Y - X'\beta)$$

This covers all of the estimators we've seen so far, where we might have $Z = X$. In data matrix notation, the estimator minimizes

$$(Z'Y - Z'X\beta)'W(Z'Y - Z'X\beta)$$

so the first order condition is:

$$\begin{aligned}
-2X'ZW(Z'Y - Z'X\hat{\beta}) &= 0 \\
\implies X'ZWZ'X\hat{\beta} &= X'ZWZ'Y \\
\implies \hat{\beta} &= (X'ZWZ'X)^{-1}X'ZWZ'Y \\
&= (S'_{XZ}WS_{XZ})^{-1}S'_{XZ}Ws_{XY}
\end{aligned}$$

where the last line is the same in Hayashi's notation. More instructively, if we set $\mu = Z'Y$ and $G = Z'X$, we want to minimize

$$(\mu - G\beta)'W(\mu - G\beta)$$

which looks exactly like weighted least squares with k regressors, ℓ observations, and weights W , so that

$$\hat{\beta} = (G'WG)^{-1}G'W\mu$$

We can additionally directly compare

$$\begin{aligned}
\hat{\beta}_{GMM}(W) &= (X'ZWZ'X)^{-1}X'ZWZ'Y \\
\hat{\beta}_{TSLS} &= (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y
\end{aligned}$$

so the TSLS estimator is GMM with weights $(Z'Z)^{-1}$. This raises two immediate questions: Since we have a family of weight matrices, which one is the best? When, if ever, is that best weight $(Z'Z)^{-1}$?

We can think about the WLS example – what are the best weights in WLS? With i.i.d. data⁶ they are the inverse standard deviation. More generally, the ideal W is the variance-covariance matrix of errors ε .

This raises the idea of estimating the optimal W . But didn't we have earlier than generalized least squares is rarely used in practice? Yes! But this example has limitations. The variance-covariance matrix of ε is $n \times n$, and estimating it comes with a number of additional assumptions (and/or nonparametric convergence rates). In our WLS example, the error is $\eta \equiv \mu - G\beta$, where the variance-covariance matrix is $\ell \times \ell$. Estimating that is plausible.

This result motivates *two-state (efficient) GMM*:

1. Compute a preliminary estimate of β to get residuals
2. Use residuals to estimate the optimal weight matrix \hat{W}
3. Report a final estimate $\hat{\beta}_{GMM}(\hat{W})$

We will show that this procedure minimizes asymptotic variance, and if we assume homoskedasticity the TSLS weights $(Z'Z)^{-1}$ will indeed be optimal.

Assumption 4.1. *We have the following assumptions for GMM:*

1. *We observe i.i.d. realizations $(Y_i, X_i, Z_i), i = 1, 2, \dots$*
2. $\mathbb{E}(Z(Y - X'\beta)) = 0$
3. $\mathbb{E}|Y|^4 < \infty, \mathbb{E}\|X\|^4 < \infty, \mathbb{E}\|Z\|^4 < \infty$
4. $Q \equiv \mathbb{E}(ZX')$ has full rank k
5. W is positive definite
6. $\Omega \equiv \mathbb{E}(ZZ'\varepsilon^2)$ is positive definite.

Remark. Assumptions 4.1 are also the assumptions for TSLS, which will emerge as a special case.

⁶Note: excludes panel estimators!

The GMM estimator has the following asymptotic distribution (a generalization of the algebra from OLS):

$$\begin{aligned}
\hat{\beta}_{GMM}(W) &= (X'ZWZ'X)^{-1}X'ZWZ'Y \\
&= \beta + (X'ZWZ'X)^{-1}X'ZWZ'\varepsilon \\
&= \beta + \left(\frac{1}{n}X'ZW\frac{1}{n}Z'X\right)^{-1}\frac{1}{n}X'ZW\frac{1}{n}Z'\varepsilon \\
&= \beta + (\mathbb{E}(XZ')W\mathbb{E}(ZX'))^{-1}\mathbb{E}(XZ')W\frac{1}{n}Z'\varepsilon + o_P(1) \\
&= \beta + (Q'WQ)^{-1}Q'W\frac{1}{n}Z'\varepsilon + o_P(1) \\
&\implies \hat{\beta}_{GMM}(W) - \beta \xrightarrow{P} 0 \\
\sqrt{n}(\hat{\beta}_{GMM}(W) - \beta) &\xrightarrow{d} \mathcal{N}(0, (Q'WQ)^{-1}Q'W\Omega WQ(Q'WQ)^{-1})
\end{aligned}$$

Remark. This only requires second moments, not fourth moments.

The matrix $\Omega = \mathbb{E}(ZZ'\varepsilon^2)$ is really the variance-covariance matrix of the moment conditions. Intuitively, a condition under which Ω has larger variance across the diagonal is noisier. In fact, in the WLS example from before, Ω parameterizes the heteroskedasticity in our regression with ℓ observations and k parameters. This suggests Ω^{-1} as the efficient weighting matrix – which we do not know, but can estimate using residuals from a preliminary regression. Furthermore, the earlier results hold if $\hat{W} \xrightarrow{P} W$.

Theorem 4.1. *Let Assumptions 4.1 hold, and let $\hat{W} \xrightarrow{P} W^* \equiv \Omega^{-1}$. Then:*

1. *The asymptotic variance becomes*

$$V^* = (Q'W^*Q)^{-1}Q'W^*\Omega W^*Q(Q'W^*Q)^{-1}$$

which simplifies to

$$V^* = (Q'\Omega^{-1}Q)^{-1}$$

2. *V^* as defined in (1) is the best asymptotic variance: $V \geq V^*$ for any other estimator (meaning that $V - V^*$ is positive semi-definite)⁷*
3. *V^* is only attained by estimators that are asymptotically equivalent to $\hat{\beta}(\Omega^{-1})$.*

Proof. We will show that (i) with efficient weighting, V simplifies to V^* , (ii) $V \geq V^*$, and (iii) the inequality is strict unless the estimators are (asymptotically) equivalent. First, write:

$$\begin{aligned}
V &= A'\Omega A, \text{ where } A = WQ(Q'WQ)^{-1} \\
V^* &= B'\Omega B, \text{ where } B = \Omega^{-1}Q(Q'\Omega^{-1}Q)^{-1}
\end{aligned}$$

and observe that

$$B'\Omega A = (Q'\Omega^{-1}Q)^{-1}Q'\Omega^{-1}\Omega WQ(Q'WQ)^{-1} = V^* = B'\Omega B \implies B'\Omega(A - B) = 0$$

Thus,

$$V = A'\Omega A = (B + (A - B))'\Omega(B + (A - B)) = \underbrace{B'\Omega B}_{V^*} + \underbrace{(A - B)'\Omega B}_0 + \underbrace{B'\Omega(A - B)}_0 + \underbrace{(A - B)'\Omega(A - B)}_{\text{p.s.d.}}$$

□

⁷This ordering is very strong! The efficient estimator will also minimize the variance of any $p'\theta$, through linearity of the quadratic form.

Remark. This motivates the efficient (two-stage) GMM estimator:

$$\begin{aligned}\hat{\beta}_{TSGMM} &\equiv \hat{\beta}(\hat{W}) \\ \hat{W} &\equiv (\mathbb{E}_n(ZZ'\hat{\varepsilon}^2))^{-1} \\ \hat{\varepsilon} &= Y - X\hat{\beta}\end{aligned}$$

where $\hat{\beta}$ is any consistent estimator of β , for example a GMM estimator with any reasonable weighting matrix. The industry standard is TSLS.

Remark. We could also use the *centered estimator*

$$\hat{W} = [\mathbb{E}_n((Z\hat{\varepsilon} - \mathbb{E}_n(Z\hat{\varepsilon}))(Z\hat{\varepsilon} - \mathbb{E}_n(Z\hat{\varepsilon}))')]^{-1}$$

which literally estimates the variance, rather than the uncentered second moment of Z . The two are the same if $\mathbb{E}Z\varepsilon = 0$, but the centered estimator is consistent for variance even if the model is misspecified.

Remark. If we further assume homoskedasticity, so that

$$\mathbb{E}(\varepsilon^2 | Z) = \sigma^2 \implies \Omega = \mathbb{E}(ZZ'\varepsilon^2) = \sigma^2 \mathbb{E}(ZZ')$$

the estimator with the ideal weighting matrix simplifies:

$$\begin{aligned}\hat{\beta}_{GMM}(\Omega^{-1}) &= (X'Z\sigma^{-2}(\mathbb{E}ZZ')^{-1}Z'X)^{-1}X'Z\sigma^{-2}(\mathbb{E}ZZ')^{-1}Z'Y \\ &= (X'Z(\mathbb{E}ZZ')^{-1}Z'X)^{-1}X'Z(\mathbb{E}ZZ')^{-1}Z'Y\end{aligned}$$

but $\mathbb{E}ZZ'$ can be estimated by $\mathbb{E}_nZZ' = \frac{1}{n}Z'Z$. Because $\frac{1}{n}$ cancels from the expression, we can succinctly write

$$\hat{\beta}_{GMM}(\Omega^{-1}) = (X'Z(Z'Z)^{-1}Z'X)^{-1}X'Z(Z'Z)^{-1}Z'Y = \hat{\beta}_{TSLS}$$

So under homoskedasticity, the efficient estimator is the two-stage least squares estimator!

Question. Should we *always* do Efficient GMM?

The case for efficient GMM is more compelling than for FGLS, and some results (notably specification testing) require efficient GMM. In practice, efficient GMM is quite common. However, there are some caveats. In finite sample, estimating Ω introduces noise. Monte Carlo simulations suggest that with small samples and moderate heteroskedasticity, TSLS may perform better. On the other hand, estimation of Ω could in principle be iterated, which (under conditions that guarantee convergence at least asymptotically) removes path dependency. Implementations for this exist, but asymptotic analysis does not suggest a gain and we might be concerned about error propagation.

Definition. With modern computing power, we could also directly compute

$$\hat{\beta}_{cGMM} \equiv \underset{\beta}{\operatorname{argmin}} \bar{g}_n(\beta)' \hat{W} \bar{g}_n(\beta)$$

where we optimize over the criterion function and the weighting matrix at the same time. This is called *Continuous Updating GMM*, and while not numerically the same estimator has the same asymptotic distribution.

Question. Should we use all the instruments we can think of?

Now that we can use more instruments than regressors, we might be tempted to use all instruments we can think of, including ‘technical instruments’ (polynomials of instruments, etc). Intuitively, we shouldn’t do that. But what is the formal argument against? First, every instrument must be justified, so the ‘cost of assumptions’ will increase (especially the exclusion restriction). Second, estimating another weight matrix introduces more finite sample noise. Third, even without weighting it can be shown that TSLS is inconsistent

⁸Unique result under conditions, which we just assume here.

if there are many instruments in the sense that $\ell_n/n \rightarrow \alpha > 0$. For this, the instruments don't even need to be weak. Actually, in practice we should only invoke asymptotics if $n \gg \ell$, limiting the number of instruments we can use.

Definition. What is qualitatively true in overidentified ($\ell > k$) models is that the model itself can be tested! We can actually introduce the *joint validity of moments test*. This works because we can (empirically) test the assumption that the instruments are orthogonal to the sample.

Theorem 4.2. *Under Assumptions 4.1, assuming the model is overspecified, then*

$$J_n \equiv J(\hat{\beta}_{TSGMM}) \xrightarrow{d} \chi_{\ell-k}^2$$

Remark. The intuition here is that we try and set an ℓ -vector to zero but only have k free parameters to do so. This means we have a residual with $\ell - k$ degrees of freedom. If the model is well specified, the residual is of order $O(n^{-1/2})$. If (and only if!) we use the efficient weighting matrix, it is further asymptotically multivariate standard normal in a certain $(\ell - k)$ -subspace. Then its square is $\chi_{\ell-k}^2$.

Proof. Previous results imply that $\frac{1}{n}Z'\varepsilon = O_P(n^{-1/2})$, so we write⁹

$$\begin{aligned} J_n &= n \left(\frac{1}{n} Z' \hat{\varepsilon} \right)' \hat{\Omega}^{-1} \left(\frac{1}{n} Z' \hat{\varepsilon} \right) \\ &\approx n \left(\frac{1}{n} Z' \hat{\varepsilon} \right)' \Omega^{-1} \left(\frac{1}{n} Z' \hat{\varepsilon} \right) \\ &= n \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right)' (C' \Omega C)^{-1} \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right) \\ &= n \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right)' \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right) \end{aligned}$$

where $\Omega^{-1} = CC' \iff \Omega(C')^{-1}C^{-1}$, meaning that C is the *Cholesky Root* of Ω^{-1} . Next, we have that

$$\begin{aligned} \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right) &= C' \frac{1}{n} Z' (\varepsilon - X(\hat{\beta} - \beta)) \\ &= C' \frac{1}{n} Z' \left[\varepsilon - X \left(\left(\frac{1}{n} X' Z \right) \hat{\Omega}^{-1} \left(\frac{1}{n} Z' X \right) \right)^{-1} \left(\left(\frac{1}{n} X' Z \right) \hat{\Omega}^{-1} \left(\frac{1}{n} Z' \varepsilon \right) \right) \right] \\ &= \left[I_\ell - C' \left(\frac{1}{n} Z' X \right) \left[\left(\left(\frac{1}{n} X' Z \right) \hat{\Omega}^{-1} \left(\frac{1}{n} Z' X \right) \right)^{-1} \left(\left(\frac{1}{n} X' Z \right) \hat{\Omega}^{-1} \left(\frac{1}{n} Z' C \right) \right) \right] \right] C' \frac{1}{n} Z' \varepsilon \\ &\approx \left[I_\ell - \underbrace{C' \left(\frac{1}{n} Z' X \right)}_{:= \hat{R}} \left[\left(\left(\frac{1}{n} X' Z \right) CC' \left(\frac{1}{n} Z' X \right) \right)^{-1} \left(\left(\frac{1}{n} X' Z \right) CC' \left(\frac{1}{n} Z' C \right) \right) \right] \right] C' \frac{1}{n} Z' \varepsilon \\ &= \left(I_\ell - \hat{R}(\hat{R}'\hat{R})^{-1}\hat{R}' \right) C' \frac{1}{n} Z' \varepsilon \\ &\approx \left(I_\ell - R(R'R)^{-1}R' \right) C' \frac{1}{n} Z' \varepsilon, \text{ where } R \equiv C'\mathbb{E}(ZX') \end{aligned}$$

We so far have

$$J_n \approx n \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right)' \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right) \quad \text{and} \quad \left(C' \frac{1}{n} Z' \hat{\varepsilon} \right) \approx \left(I_\ell - R(R'R)^{-1}R' \right) C' \frac{1}{n} Z' \varepsilon$$

⁹In this proof, \approx means that we drop an $o_P(1)$ term.

and observe that

$$\sqrt{n}C' \left(\frac{1}{n} Z' \varepsilon \right) \xrightarrow{d} \mathcal{N}(0, C' \Omega C) = \mathcal{N}(0, C' (C')^{-1} C^{-1} C) = \mathcal{N}(0, I_\ell)$$

Define the random variable $u \sim \mathcal{N}(0, I_\ell)$, then it follows that

$$J_n \xrightarrow{d} (I_\ell - R(R'R)^{-1}R')' (I_\ell - R(R'R)^{-1}R') u \sim \chi_{\ell-k}^2$$

because $I_\ell - R(R'R)^{-1}R'$ projects u onto a lower dimensional subspace. It remains to show that the subspace is of dimension $\ell - k$. This follows from the fact that $I_\ell - R(R'R)^{-1}R'$ is the annihilator matrix associated with $R = C'\mathbb{E}(ZX')$, which has rank k and null space of rank $\ell - k$. (You could use the rank-nullity condition of idempotent matrices to prove this statement directly). \square

Example. Visualization of the above theorem. Consider a situation with $k = 1$ endogenous regressor and $\ell = 2$ instruments, where:

$$\Omega = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow \Omega^{-1} = \begin{bmatrix} 1/4 & 0 \\ 0 & 1 \end{bmatrix}, C = \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mathbb{E}ZX' = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow R = \begin{bmatrix} 1/2 \\ 1 \end{bmatrix}$$

The first moment condition has four times the variance of the second, so in the WLS analogy it should be given half of the weight. The column space of R is the line spanned by $[1/2 \ 1]'$, so its null space is the line spanned by $[-1 \ 1/2]'$. Visually, we have Figure 5.

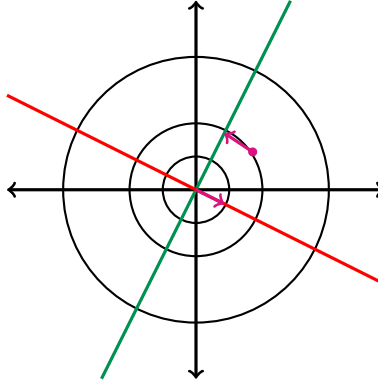


Figure 5: **Column Space** and **Null Space** of R in a simple example. An **element** is projected onto the column space, leaving residuals in the null space that are also standard normal.

Definition. Let $\theta = r(\beta)$ and $\hat{\theta} = r(\hat{\beta})$ for some function $r : \mathbb{R}^k \rightarrow \mathbb{R}^q$, where $r \in \mathbf{C}^1$. Suppose also that $R(\beta) = \frac{\partial r(\beta)}{\partial \beta'}$ has full rank of q at the true value of β . Then we can construct a **Wald test**:

$$W \equiv n(\hat{\theta} - \theta) \left(R(\hat{\beta})' \hat{V}_\beta R(\hat{\beta}) \right)^{-1} (\hat{\theta} - \theta) \xrightarrow{d} \chi_q^2$$

We will later develop other tests in a more general setting.

Remark. We can extend this analysis to multiple equations, where

$$\begin{aligned} Y_m &= X_m \beta_m + \varepsilon_m & m &= 1, \dots, M \\ \mathbb{E}Z_m \varepsilon_m &= 0 & m &= 1, \dots, M \end{aligned}$$

The X_m may be the same (like in seemingly unrelated regression equations), or they may overlap (like in panel data with some time invariant regressors) – coefficients are not restricted across equations in a general setting, but they are in many applications.

Example. In *Seemingly Unrelated Regressions* (this example from [Griliches, 1976](#)), we may have that

$$\begin{aligned} LW69 &= \alpha_1 + \beta_1 \cdot schooling69 + \gamma_1 \cdot IQ + \delta_1 \cdot experience69 + \varepsilon_1 \\ KWW &= \alpha_2 + \beta_2 \cdot schooling69 + \gamma_2 \cdot IQ + \varepsilon_2 \end{aligned}$$

where $LW69$ is the logged wage, and KWW is a measure of ability (as are IQ and $experience69$). We can think of this as regressing the two outcomes on the same regressor but assuming *a priori* that $\delta_1 = 0$. That alone makes this model overidentified, and jointly estimating both equations leads to a more efficient estimate (*if we believe the assumptions!*). We can construct a (fictitious) version of this with *panel data*, and get

$$\begin{aligned} LW69 &= \alpha_1 + \beta_1 \cdot schooling69 + \gamma_1 \cdot IQ + \delta_1 \cdot experience69 + \varepsilon_1 \\ LW80 &= \alpha_2 + \beta_2 \cdot schooling80 + \gamma_2 \cdot IQ + \delta_2 \cdot experience80 + \varepsilon_2 \end{aligned}$$

We can define

$$\begin{aligned} \bar{Y} &= \begin{bmatrix} Y_1 \\ \vdots \\ Y_M \end{bmatrix} & \bar{X} &= \begin{bmatrix} X_1 & 0 & \cdots & 0 \\ 0 & X_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & X_M \end{bmatrix} \\ \bar{\beta} &= \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_M \end{bmatrix} & \bar{Z} &= \begin{bmatrix} Z_1 & 0 & \cdots & 0 \\ 0 & Z_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Z_M \end{bmatrix} \end{aligned}$$

and we have the moment condition

$$\mathbb{E}(\bar{Z}(\bar{Y} - \bar{X}'\bar{\beta})) = 0$$

and estimator

$$\hat{\beta}(W) = (\mathbb{E}_n(\bar{X}\bar{Z}')W\mathbb{E}_n(\bar{Z}\bar{X}'))^{-1} (\mathbb{E}_n(\bar{X}\bar{Z}')W\mathbb{E}_n(\bar{Z}\bar{Y}))$$

The estimator is the same as before, except that the data matrix notation becomes unwieldy. Assuming that $(\bar{Y}, \bar{X}, \bar{Z})$ fulfill Assumptions 4.1, we have that

$$\begin{aligned} \sqrt{n}(\hat{\beta}(W) - \beta) &\xrightarrow{d} \mathcal{N}(0, V_\beta) \\ V_\beta &= (\bar{Q}'W\bar{Q})^{-1}\bar{Q}'W\bar{\Omega}W\bar{Q}(\bar{Q}'W\bar{Q})^{-1} \\ \bar{Q} &= \mathbb{E}(\bar{Z}\bar{X}') \\ \bar{\Omega} &= \mathbb{E}(\bar{Z}\varepsilon\varepsilon'\bar{Z}') \end{aligned}$$

and results on efficient GMM are also the same as before.

Remark. This is extremely powerful! We just derived a collection of historically distinct estimators. A small caveat is that we must think very carefully about what the assumptions on the new objects actually mean. We've made the assumption that $(X_1, \dots, X_M, Y_1, \dots, Y_M, Z_1, \dots, Z_M)$ are *all* i.i.d., which is stronger than assuming that equation-by-equation.

Remark. This is the same as estimating the equations separately (*i.e.* will lead to the same result) if (i) everything is just identified *or* (ii) W is block diagonal, where its blocks correspond to equations. It can be

instructive to write out the estimator for $M = 2$:

$$\begin{aligned}\hat{\beta}(\hat{W}) &= \left[\begin{bmatrix} \mathbb{E}_n(X_1 Z_1') & 0 \\ 0 & \mathbb{E}_n(X_2 Z_2') \end{bmatrix} \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') & 0 \\ 0 & \mathbb{E}_n(X_2 Z_2') \end{bmatrix} \right]^{-1} \\ &\quad \cdot \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') & 0 \\ 0 & \mathbb{E}_n(X_2 Z_2') \end{bmatrix} \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} \begin{bmatrix} \mathbb{E}_n(Z_1 Y_1) \\ \mathbb{E}_n(Z_2 Y_2) \end{bmatrix} \\ &= \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 X_1') & \mathbb{E}_n(X_1 Z_1') W_{12} \mathbb{E}_n(Z_2 X_2') \\ \mathbb{E}_n(X_2 Z_2') W_{21} \mathbb{E}_n(Z_1 X_1') & \mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 X_2') \end{bmatrix} \cdot \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 Y_1) + \mathbb{E}_n(X_1 Z_1') W_{12} \mathbb{E}_n(Z_2 Y_2) \\ \mathbb{E}_n(X_2 Z_2') W_{21} \mathbb{E}_n(Z_1 Y_1) + \mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 Y_2) \end{bmatrix}\end{aligned}$$

What happens if $W_{12} = W_{21} = 0$? We get the simplification:

$$\begin{aligned}&= \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 X_1') & 0 \\ 0 & \mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 X_2') \end{bmatrix} \cdot \begin{bmatrix} \mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 Y_1) \\ \mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 Y_2) \end{bmatrix} \\ &= \begin{bmatrix} (\mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 X_1'))^{-1} \mathbb{E}_n(X_1 Z_1') W_{11} \mathbb{E}_n(Z_1 Y_1) \\ (\mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 X_2'))^{-1} \mathbb{E}_n(X_2 Z_2') W_{22} \mathbb{E}_n(Z_2 Y_2) \end{bmatrix} \\ &= \begin{bmatrix} \hat{\beta}_1(W_{11}) \\ \hat{\beta}_2(W_{22}) \end{bmatrix}\end{aligned}$$

So when we assume that the cross-equation weights are zero, the multi-equation GMM just stacks the single equation estimators!

Question. When should we estimate equations jointly?

Answer. If we naively interpret this result, we could estimate all of the equations in the world jointly – if they're actually unrelated, the weighting matrix will pick that up. This seems intuitively untrue, but how do we show it formally? Well, we face an escalating number of nuisance parameters (entries of W). More importantly, model misspecification is contagious – the estimator's probability limit equals

$$\text{plim } \hat{\beta}(W) = \beta + (\mathbb{E}(\bar{X} \bar{Z}') W \mathbb{E}(\bar{Z} \bar{X}'))^{-1} \mathbb{E}(\bar{X} \bar{Z}') W \mathbb{E}(\bar{Z} \varepsilon)$$

If any one entry of $\mathbb{E}(\bar{Z} \varepsilon)$ is nonzero, then every entry of the matrix product is nonzero. This holds unless W is block diagonal in the equations. In that case (and *only* that case!) the joint estimation is efficient.

Example. *Multiple Regression vs. Seemingly Unrelated Regression* If $X_1 = \dots = X_M = Z_1 = \dots = Z_M$, then this is just multiple regression – that is, we regress different Y_1, \dots, Y_M on the same exogenous regressors. We can verify that the estimator just stacks OLS in this case.

Alternatively, suppose that some regressors are dropped from some equations, meaning that we think their coefficients are zero. However, we still consider them exogenous in all equations, so we can use them as overidentifying instruments. Formally, let $Z_1 = \dots = Z_M = \bigcup_{m=1}^M X_m$ and the X_m are not all the same. This is the *Seemingly Unrelated Regression (SUR)* estimator.¹⁰

Remark. As we add assumptions to the basic model, we can recover some classic estimators that were developed independently. These assumptions typically allow us to simplify expressions. See Hayashi for more detail.

- If we assume homoskedasticity, we get *Full Information Instrumental Variables Efficient (FIVE) estimation* (from Brundy & Jorgenson, 1971)
- If we additionally assume that $Z_1 = \dots = Z_M$, we get *Three-Stage Least Squares (3SLS)*¹¹ (from Zellner & Theil, 1962)
- SUR is the final specialization, where we set $Z_1 = \dots = Z_M = \bigcup_{m=1}^M X_m$

¹⁰Historically, it would also require homoskedasticity, but that's actually a distinct issue.

¹¹In modern terminology, this is a two-stage estimator. However, it could be expressed as TSLS with 'pre-pre-estimation' of cross-equation correlation of errors, hence the name.

Remark. Next, we can think about common coefficients. Consider the following specification:

$$Y_m = X_m\beta + \varepsilon_m, \quad \mathbb{E}Z_m\varepsilon_m = 0 \quad \forall m = 1, \dots, M$$

The main intuition here is that an essentially the same (but not constant) covariate is observed across equations. We do not need to revisit every covariate in every equation: some components of X_m could be zero almost surely. If we define

$$\bar{Y} = \begin{bmatrix} Y_1 \\ \vdots \\ Y_M \end{bmatrix} \quad ; \quad \bar{X} = [X_1 \quad \cdots \quad X_M] \quad ; \quad \bar{Z} = \begin{bmatrix} Z_1 & 0 & \cdots & 0 \\ 0 & Z_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & Z_M \end{bmatrix}$$

then we have moment condition $\mathbb{E}(\bar{Z}(\bar{Y} - \bar{X}'\beta)) = 0$ and estimator

$$\hat{\beta}(W) = (\mathbb{E}_n(\bar{X}\bar{Z}')W\mathbb{E}_n(\bar{Z}\bar{X}'))^{-1} \mathbb{E}_n(\bar{X}\bar{Z}')W\mathbb{E}_n(\bar{Z}\bar{Y})$$

This looks like what we had before, but the different definition of \bar{X} changes identification: that $\mathbb{E}(\bar{Z}\bar{X}')$ has full rank is implied by $\mathbb{E}(Z_m X_m')$ having full rank for each $m = 1, \dots, M$.

Common coefficients allow for estimation of many parameters that would not otherwise be identifiable. An important application is panel data:

- If we impose the assumptions that characterized SUR before, we get the *Random Effects* estimator
- If we furthermore assume that ε_m is uncorrelated across m , then this simplifies to *Pooled OLS*

The difference between those estimators is about efficiency, not identification. We will revisit random effects and pooled OLS soon.

Remark. Efficient GMM vaguely resembles WLS or *Feasible Generalized Least Squares (FGLS)*. As a reminder, weighted least squares minimizes

$$(Y - X\beta)'W(W - X\beta)$$

where the weighting matrix W gives differential weights to different observations. If we know that $\mathbb{E}(\varepsilon\varepsilon' | X) = \sigma^2 \cdot \Omega$ for known Ω , then setting $W = \Omega^{-1}$ is variance minimizing, and the resulting estimator is the BLUE by Gauss-Markov. To see the analogy to WLS as undergraduates learn it, note that we can equivalently minimize

$$(C(Y - X\beta))'(C(Y - X\beta))$$

where C is the Cholesky root of W . In particular, if observations are uncorrelated,

$$\Omega = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix} \quad ; \quad \Omega^{-1} = \begin{bmatrix} 1/\sigma_1^2 & 0 & \cdots & 0 \\ 0 & 1/\sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_n^2 \end{bmatrix} \quad ; \quad C = \begin{bmatrix} 1/\sigma_1 & 0 & \cdots & 0 \\ 0 & 1/\sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1/\sigma_n \end{bmatrix}$$

and the minimand can be expressed as $\sum_{i=1}^n ((y_i - x_i'\beta)/\sigma_i)^2$.

Question. What do we do if we don't know Ω ? We could estimate it, but estimating an $n \times n$ matrix from n data points seems absurd.

There's two possible paths: (i) if we have a specific parameter model for how $\mathbb{E}(\varepsilon^2 | X = x)$ changes with x , we could estimate it; and (ii) we could run a general flexible regression of $\hat{\mu}_i^2$ on X . The latter regressions are often run to construct Breusch-Pagan or other heteroskedasticity tests. However, (i) is rare and for (ii) we incur a lot of effort (and error propagation!) for reweighting that only matters if weights are very different

than even. This explains why the latter approach has a name (Feasible Generalized Least Squares) but is extremely uncommon in practice.

So what does this have to do with GMM? We can write the GMM estimator as the weighted OLS estimator in a fictitious regression, where we set $\mu = Z'Y$ and $G = Z'X$, then we want to minimize

$$(\mu - G\beta)'W(\mu - G\beta)$$

and this looks like WLS with k regressors, ℓ observations, and weight matrix W . The closed-form WLS estimator

$$\hat{\beta}_{WLS} = (G'WG)^{-1}G'W\mu$$

precisely recovers the GMM estimator (once we substitute back the transformations). So if GMM is similar to FGLS, why is it so much more common? The analogy is useful but has limitations. Every moment condition in the original problem is an ‘observation’ in the fictitious regression. In our WLS analogy, the ‘error’ is $\eta \equiv \mu - G\beta$, with variance-covariance matrix of size $\ell \times \ell$, so we are estimating a much smaller matrix, which becomes a lot more precise as n grows! We can also see this in the objective functions:

$$(Y - X\beta)' \underbrace{W}_{n \times n} (Y - X\beta) \quad \text{vs.} \quad (\mu - G\beta)' \underbrace{W}_{\ell \times \ell} (\mu - G\beta)$$

5 Panel Data

Remark. An *extremely quick* and non-exhaustive treatment.

Definition. *Panel data* are data that come in a two-dimensional array, most commonly in ‘time’ and ‘units of observation’. We have *sample size* n and *number of waves* T . With a *balanced panel*, we have nT observations.¹² We call a *short panel* one where in practice $T \ll n$, and the asymptotics require that we fix T and let $n \rightarrow \infty$. Alternatively, we have a *long panel*, where we fix n and let $T \rightarrow \infty$. This is really a multivariate time series. Some applications require analyses where we send $n \rightarrow \infty$ and $T \rightarrow \infty$ at the same or different rates. Here, we consider only short panels.

Remark. All estimators in this section are variations on multiple equation common coefficients GMM. We will develop from specific to general here, starting with

$$Y_{it} = X'_{it}\beta + \varepsilon_{it} \quad \text{where } \mathbb{E}(X_{it}\varepsilon_{it}) = 0$$

From there, we get:

Definition. The *pooled OLS estimator* is defined as

$$\hat{\beta}_{pool} \equiv \left(\sum_{i=1}^n X'_i X_i \right)^{-1} \sum_{i=1}^n X'_i Y_i = (X'X)^{-1} X'Y \xrightarrow{p} \beta$$

This is unbiased and consistent under the assumptions that confirm that previously (where X and Y are pooled data matrices). It also has the same asymptotics if we make the same assumptions as above. However, the homoskedasticity assumption is *incredibly* strong here. At the very least, we need to use a *cluster-robust variance estimator*

$$\hat{V}_{pool} = (X'X)^{-1} \left(\sum_{i=1}^n X'_i \hat{\varepsilon}_i \hat{\varepsilon}_i X_i \right) (X'X)^{-1}$$

We had previously omitted cluster-dependent errors, but note the analogy to FGLS: the central term estimates the $(T \times T)$ -matrix $\mathbb{E}(X'_i \hat{\varepsilon}_i \hat{\varepsilon}_i X_i)$. However, we could do even better! As previously, the variance matrix is ‘small’, which raises the possibility of using an FGLS-like approach for estimation. In fact, we are precisely doing efficient (multi-equation) GMM.

Definition. A popular model is the *Random Effects Model*, where we assume cross-sectional homoskedasticity but with the specific structure

$$\mathbb{E}(\varepsilon_i \varepsilon'_i | X_i) = \begin{bmatrix} \sigma_u^2 + \sigma_\epsilon^2 & \sigma_u^2 & \cdots & \sigma_u^2 \\ \sigma_u^2 & \sigma_u^2 + \sigma_\epsilon^2 & \cdots & \sigma_u^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_u^2 & \sigma_u^2 & \cdots & \sigma_u^2 + \sigma_\epsilon^2 \end{bmatrix}$$

which has only two degrees of freedom. The assumption here is that $\varepsilon_{it} = u_i + \epsilon_{it}$ where u_i and ϵ_{it} are uncorrelated. In other words, our structural equation becomes the model

$$Y_{it} = X'_{it}\beta + u_i + \epsilon_{it} \quad \text{where } \mathbb{E}(X_{it}u_i) = \mathbb{E}(X_{it}\epsilon_{it}) = 0$$

Remark. The random effects estimator can be thought of as the feasible GLS estimator for this model, where we pre-estimate σ_u and σ_ϵ . We will get more details later, under some theory yet to come. Note that Hayashi defines the random effects estimator as a two-stage GMM assuming homoskedasticity, which is similar to three-stage least squares. This treatment has more degrees of freedom.

What’s important is that (i) random effects has the same identifying assumptions as pooled OLS, and (ii) it adds a FGLS/TSGMM step, and since within-unit correlation of ε is both salient and easily modeled (note:

¹²We will consider only balanced panels.

in short panels only!), this is often desirable.

Remark. In the random effects model, the condition that $\mathbb{E}X_{it}u_i = 0$ is very restrictive. If we think of u_i as unobserved, time-invariant covariates, we are saying that these cannot at all be correlated with the observed. If this fails, do we have anything else to estimate?

Consider the *Fixed Effects (FE)* equation

$$Y_{it} - Y_{i,t-1} = (X_{it} - X_{i,t-1})'\beta + \epsilon_{it} - \epsilon_{i,t-1}$$

When can we estimate this by OLS? We need that (i) ϵ_{it} is uncorrelated with past and future ϵ_t , and (ii) that $(X_{it} - X_{i,t-1})$ fulfills a rank condition (this will fail with time-invariant regressors). This is the basic idea of fixed effects estimation, where we can ‘difference away’ the fixed effects in one of three ways: (i) first differencing (*between estimator*), (ii) demeaning (*within estimator*), and (iii) adding an indicator of each unit (*dummy variable regression*). Demeaning and dummy variables are numerically the same, and correspond to the ‘classic’ fixed effects estimator. From the point of view of identification, the methods are the same, but they imply different weighting matrices. If the weighting matrix is pre-estimated, they are asymptotically the same.

We additionally have that:

1. The implied weighting matrix $M = 1(1'1)^{-1}1'$ is efficient if the idiosyncratic error ϵ_{it} is homoskedastic and uncorrelated. (for the between estimator, we need that ϵ_{it} follows a random walk in direction t . This is less salient)
2. We can actually show that FE equals TSLS (or really SUR), thinking of the cross-equation restrictions in

$$\mathbb{E}X_{is}\epsilon_{it} = 0 \quad \forall s, t$$

as overidentifying restrictions / instruments

3. The estimator necessarily has higher (asymptotic) variance than pooled OLS. This is because OLS algebra applies, but demeaning reduces the sum of squared deviations of any random variable
4. For variance estimation, a degrees of freedom adjustment of $T(T-1)$ is not negligible for realistic T and is therefore recommended for including under asymptotic justification
5. In the formal model that motivates the random effects model, the FE estimator can be used to estimate σ_ϵ^2 . Doing this first and then backing out σ_u^2 is the standard approach

Assumption 5.1. *We assume the following:*

1. $Y_{it} = X'_{it}u_i + \epsilon_{it}$ for all $t = 1, \dots, T$, $T \geq 2$
2. X_{it} are i.i.d.
3. $\mathbb{E}(X_{is}\epsilon_{it}) = 0$ for all $s, t = 1, \dots, T$
4. $Q \equiv \mathbb{E}(\bar{X}'_i \bar{X}_i)$ is positive definite, where $\bar{X}_i = (X_{i2} - \bar{X}_i, \dots, X_{iT} - \bar{X}_i)'$
5. $\mathbb{E}\epsilon_{it}^4 < \infty$, $\mathbb{E}\|X_{it}\|^4 < \infty$

Theorem 5.1. *Under Assumptions 5.1,*

$$\sqrt{n}(\hat{\beta}_{fe} - \beta) \xrightarrow{d} \mathcal{N}(0, Q^{-1}\Omega Q^{-1})$$

$$Q \equiv \mathbb{E}(\bar{X}'_i \epsilon_i \epsilon_i' \bar{X}_i)$$

Remark. Though we don’t go into it here, it is true that undetrended time series can exhibit extreme (and considerably confounded) correlations. This motivates us to detrend our panel and estimate the equation

$$Y_{it} = X_{it}\beta + u_i + \nu_t + \epsilon_{it}$$

This is called the *Two-Way Fixed Effects (TWFE)* model, and when $T = 2$ (only!) it is equivalent to the *difference-in-difference (DiD)* estimation. The best known estimator is the pooled OLS after a double-within transformation. This is the subject of an active literature (see: [Callaway & Sant'Anna \(2021\)](#), [Goodman-Bacon \(2021\)](#), etc) and you should consult with an econometrician before using it.

Definition. The *dynamic panel* is estimated by the equation

$$Y_{it} = \alpha \cdot Y_{i,t-1} + X_{it} \cdot \beta + u_i + \epsilon_{it}$$

We consider only one lag for exposition. In practice if data have a trend it would be essential to separately model that to avoid spurious correlation. The qualitatively new problem is that, even under the preceding assumptions, we have that Y_{it} is an endogenous regressor in the transformed data. It is most easily seen with first differencing:

$$\mathbb{E}(\Delta Y_{i2} \Delta \epsilon_{i3}) = \mathbb{E}((Y_{i2} - Y_{i1})(\epsilon_{i3} - \epsilon_{i2})) = \underbrace{\mathbb{E}(Y_{i2} \epsilon_{i3})}_0 - \underbrace{\mathbb{E}(Y_{i1} \epsilon_{i3})}_0 - \underbrace{\mathbb{E}(Y_{i2} \epsilon_{i2})}_{\sigma_\epsilon^2} + \underbrace{\mathbb{E}(Y_{i1} \epsilon_{i2})}_0 = -\sigma_\epsilon^2$$

This finding implies that FE estimation is biased and inconsistent. In fact, if we assume stationarity ($|\alpha| < 1$), the asymptotic bias of $\hat{\alpha}$ can be computed as

$$\text{plim } \hat{\alpha}_{\text{fe}} = \alpha - \frac{1 + \alpha}{2\alpha/(1 - \alpha) + (T - 1)/(1 - \alpha^{T-1})}$$

and a bias of the same order (*i.e.* $O(1/T)$) is inherited by $\hat{\beta}$. Note that the bias is negative and not particularly small even for moderately large T . We can conclude that FE are not appropriate in this setting. The previous is from [Nickell \(1981\)](#) and is often called the *Nickell Critique*. Some practitioners deliberately ignore this because the fixes have their own issues.

Remark. What can we do about this? Instrumental variables! We write that

$$\text{cov}(Y_{i,t-2}, \Delta Y_{i,t-1}) \neq 0 \quad ; \quad \mathbb{E}(Y_{i,t-1} \Delta \epsilon_{it}) = \mathbb{E}(Y_{i,t-2} \epsilon_{it}) - \mathbb{E}(Y_{i,t-2} \epsilon_{i,t-1}) = 0$$

So as long as $T \geq 3$, we have an instrument. (More generally, T must exceed the number of lags of Y used by at least 2). This gives rise to the *(Anderson-)Hsiao Estimator*. An important limitation is that the estimator is very sensitive to misspecification, meaning that ϵ_{it} must actually be uncorrelated and the correct number of lags must be specified. The above algebra also implies that $Y_{i,t-3}, Y_{i,t-4}, \dots$ are valid and relevant instruments (though they get weaker with distance). We can also do overidentified multiple GMM, which is the *Arellano-Bond Estimator*, which exists as a one-step (as in TSLS) or a two-step (as in Efficient GMM) estimator.

6 Extremum Estimation

Definition. An *extremum estimator* is any estimator defined as

$$\hat{\theta} = \underset{\theta \in \Theta}{\text{argmin}} Q_n(W_1, \dots, W_n; \theta)$$

for some parameter θ in parameter space Θ , where W_1, \dots, W_n is a sample. The *criterion function* $Q_n(\cdot)$ must be indexed by n because its mathematical form necessarily depends on n . However, usually it is intuitively the same function at different n . Consider for example $Q_n(\cdot) = \frac{1}{n} \sum_{i=1}^n (Y_i - X_i' b)^2$. Similarly to GMM notation, we will often drop the data and call this $Q_n(\theta)$.

Remark. The intuition for why this would estimate a true parameter value θ_0 is that:

$$\begin{array}{ccc} \theta_0 = \underset{\theta \in \Theta}{\operatorname{argmin}} Q(\theta) & & \hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} Q_n(\theta) \\ Q_n(\cdot) \rightarrow Q(\cdot) & \xRightarrow{\quad ? \quad} & \hat{\theta} \rightarrow \theta_0 \end{array}$$

That is, the *sample criterion* $Q_n(\cdot)$ estimates some *population criterion* $Q(\cdot)$ that is minimized at θ_0 . Intuitively, in ‘nice’ cases that implies that $\hat{\theta} \rightarrow \theta_0$. An illustration of a ‘nice’ case is Figure 6.

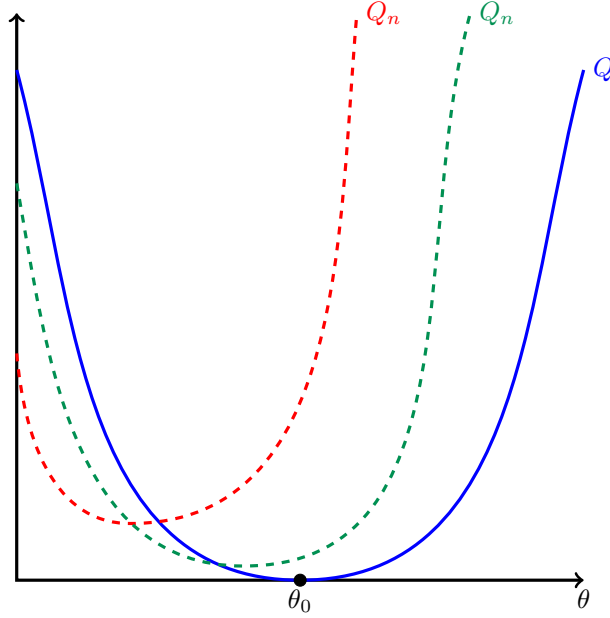


Figure 6: ‘Nice’ Extremum Estimation, with the *true criterion function*, an estimate with *small* n , and an estimate with *larger* n . The true θ_0 is additionally marked.

Example. *GMM*

$$\begin{aligned} Q(\theta) &= \mathbb{E}g(\theta)'W\mathbb{E}g(\theta) \\ Q_n(\theta) &= \mathbb{E}_n g(\theta)' \hat{W} \mathbb{E}_n g(\theta) \end{aligned}$$

Example. *Method of Simulated Moments (MSM)*

$$\begin{aligned} Q(\theta) &= (\pi(\theta) - \pi_0)'W(\pi(\theta) - \pi_0) \\ Q_n(\theta) &= (\tilde{\pi}(\theta) - \hat{\pi})' \hat{W} (\tilde{\pi}(\theta) - \hat{\pi}) \end{aligned}$$

The function $\pi(\cdot)$ maps parameter values onto implies moments of the data (*e.g.* means, variances, or even entire time series of inflation, unemployment, etc...), where π_0 are the true moments and $\hat{\pi}$ an estimate, and $\tilde{\pi}(\cdot)$ is a *simulated analog* of $\pi(\cdot)$. This (interestingly) differs from GMM if simulation noise in $\tilde{\pi}(\cdot)$ cannot be ignored. Otherwise, it is just GMM (but sometimes called MSM).

Example. *Nonlinear Least Squares*

$$Q(\theta) = \mathbb{E}(Y - m(X, \theta))^2$$

$$Q_n(\theta) = \mathbb{E}_n(Y - m(X, \theta))^2$$

You could argue that this an example of GMM (and it is!) but it was developed separately.

Example. *Maximum Likelihood*

$$Q(\theta) = \mathbb{E}\ell(W; \theta)$$

$$Q_n(\theta) = \mathbb{E}_n\ell(W; \theta)$$

The “conceptual” definition is at first glance different, but we will later derive the above from it.

Definition. The most important special case are *m-estimators*, where

$$Q(\theta) = \mathbb{E}m(W; \theta)$$

$$Q_n(\theta) = \mathbb{E}_nm(W; \theta)$$

for some known, real-valued function $m(\cdot)$. Some examples are maximum likelihood, where $m(W; \theta) = \ell(W; \theta)$, and one-step GMM where $m(W; \theta) = g(W; \theta)'Wg(W; \theta)$. Consider why efficient GMM is not an m -estimator. This class of estimators is of interest because some of the building blocks of asymptotic theory are available at exactly this level of geometry. Note that some texts (*not* Hayashi) use m -estimation as a synonym for extremum estimation.

Remark. We will next formalize the intuitive argument for consistency. We will start at a high-level and then verify in special cases. The following will assume that $\text{argmin}_{\theta \in \Theta} Q_n(\theta)$ exists, and we could see that everything goes through as long as $Q_n(\hat{\theta}) \leq \inf_{\theta \in \Theta} Q_n(\theta) + \frac{1}{n}$. Thus, $\hat{\theta}$ can be an arbitrary choice that fulfills this constraint. This settles existence and is also practically relevant because $\hat{\theta}$ may be numerically evaluated and thus not exact.

Theorem 6.1. *Consistency* Assume that:

1. The sample criterion uniformly consistently estimates the population criterion:

$$\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| \xrightarrow{P} 0$$

2. θ_0 is a unique and well-separated global minimum of $Q(\cdot)$:

$$\forall \varepsilon > 0, \exists \delta > 0 \text{ s.t. } Q^\varepsilon \equiv \inf_{\theta \in \Theta: \|\theta - \theta_0\| \geq \varepsilon} Q(\theta) \geq Q(\theta_0) + \delta$$

Then $\hat{\theta} \xrightarrow{P} \theta_0$.

Proof. Fix $\varepsilon > 0$ and define $Q_n^\varepsilon \equiv \inf_{\theta \in \Theta: \|\theta - \theta_0\| \geq \varepsilon} Q_n(\theta)$. Then we have that

$$\begin{aligned} \mathbb{P}(\|\hat{\theta} - \theta_0\| > \varepsilon) &\leq \mathbb{P}(Q_n^\varepsilon \leq Q_n(\theta_0)) \\ &= 1 - \mathbb{P}(Q_n^\varepsilon > Q_n(\theta_0)) \\ &\leq 1 - \mathbb{P}\left(Q_n^\varepsilon > Q_\varepsilon - \frac{\delta}{2}, Q_n(\theta_0) < Q(\theta_0) + \frac{\delta}{2}\right) \\ &\rightarrow 0 \end{aligned}$$

where all inequalities follow from logical implication, and the last step uses the first assumption. Thus, $\hat{\theta} \xrightarrow{P} \theta_0$. \square

Remark. The preceding result used both (i) uniform convergence and (ii) the well-separated minimum. We will provide lower-level conditions that imply these.

Theorem 6.2. *If we assume that:*

1. $Q(\cdot)$ is continuous
2. Θ is compact
3. $\theta_0 = \operatorname{argmin}_{\theta' \in \Theta} Q(\theta')$ is unique

Then θ_0 is a well-separated minimum.

Proof. Fix $\varepsilon > 0$. By the Weierstrass Theorem, Q^ε is attained by some θ^ε with $\|\theta^\varepsilon - \theta_0\| \geq \varepsilon$. Set $\delta = Q(\theta^\varepsilon) - Q(\theta_0)$, which is strictly positive by the third assumption. \square

Theorem 6.3. *If we assume that:*

1. $\hat{\theta}$ is an m -estimator
2. The data are i.i.d. realizations of W
3. $m(W; \theta)$ is almost surely continuous in θ
4. $|m(W; \theta)| \leq G(W)$ for some function G where $\mathbb{E}G(W) < \infty$
5. Θ is compact

then $Q_n(\cdot)$ converges to $Q(\cdot)$ uniformly.

Proof. This is the *Uniform Law of Large Numbers*. \square

Remark. We can consolidate the above into a single result:

Theorem 6.4. *Consolidated Consistency Assume that*

1. $Q(\cdot)$ is continuous
2. Θ is compact
3. θ_0 uniquely minimizes $Q(\theta)$
4. $\sup_{\theta \in \Theta} |Q_n(\theta) - Q(\theta)| \xrightarrow{P} 0$

Then $\hat{\theta} \xrightarrow{P} \theta_0$.

Proof. Sufficiency proved already. We can think intuitively about why each condition is necessary:

1. Continuity: Without continuity, we could have a unique but not well-separated minimum, so the function could be arbitrarily close while the estimate is far away.
2. Compactness: If the minimizer is at the boundary of an open set (*i.e.* a non-compact domain), then we will not approximate it, even with an arbitrarily good functional approximation! Alternatively, if the set is unbounded we could think of similar.
3. Uniqueness: If θ_0 is not unique, even a function that is a very good approximation of Q may be minimized close to the ‘wrong’ minimizer, as was seen when we thought about well-separated minima.
4. Uniform Convergence: Consider converging pointwise rather than uniformly. We could very possibly have Q_n limiting pointwise to a function with a true population minimum that is positive distance from any intermediate minima of the various Q_n . Think of the limiting behavior in Figure 7, where we have that the intermediate minima for each Q_n are approaching the origin, while the population minimum is strictly positive.

\square

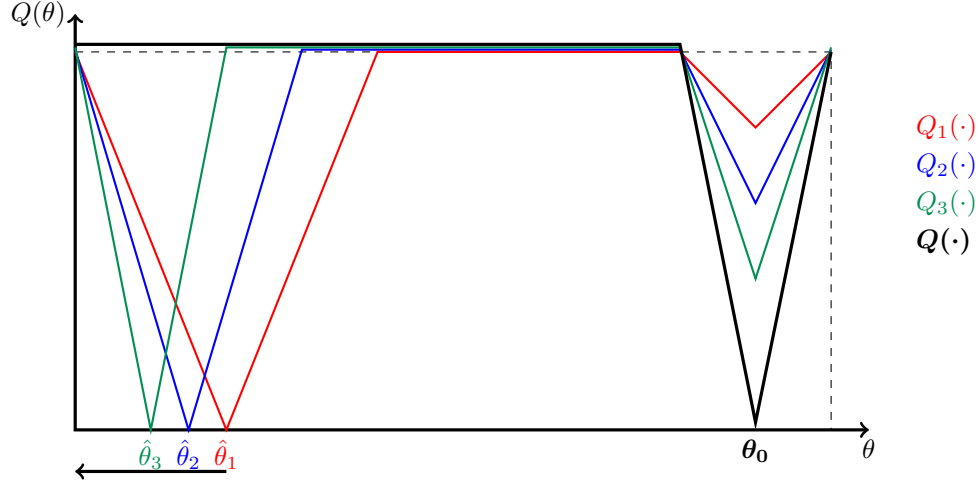


Figure 7: Pointwise (but not uniform) Convergence

Remark. When θ_0 is not unique, we could think of working with *partial identification*, which involves estimation and inference theory for $\Theta_I = \operatorname{argmin}_{\theta \in \Theta} Q(\theta)$ (a possibly non-singleton *identified set*). This is an active literature.

Theorem 6.5. *Consistency for Convex $Q(\cdot)$* If we assume that

1. Θ is convex
2. $\theta_0 \in \operatorname{int} \Theta$
3. θ_0 uniquely minimizes $Q(\theta)$
4. $Q_n(\cdot)$ is convex
5. $|Q_n(\theta) - Q(\theta)| \xrightarrow{P} 0 \forall \theta \in \Theta$

Then $\hat{\theta} \xrightarrow{P} \theta_0$.

Proof. The proof of a simplified statement is left as homework. Essentially, we need only pointwise convergence, and the assumptions guarantee that there exists $\hat{\theta}$ that exactly minimizes $Q_n(\cdot)$. \square

Remark. We are about to move on to \sqrt{n} -asymptotic normality. Are there intermediate assumptions under which we can ensure a rate of convergence without ensuring asymptotic normality? Yes! They relate the curvature of $Q(\cdot)$ at θ_0 to that rate. The rate is \sqrt{n} if $Q(\cdot)$ locally dominates some quadratic function. The proof is van der Vaart & Wellner's *Argmax Theorem* in *Weak Convergence and Empirical Processes*.

Assumption 6.1. *In order to obtain asymptotics, we need:*

1. $\hat{\theta} \xrightarrow{P} \theta_0$
2. $\theta_0 \in \operatorname{int}(\Theta)$
3. $Q_n(\cdot)$ is twice continuously differentiable in an open neighborhood N of θ_0
4. $\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \xrightarrow{d} \mathcal{N}(0, \Sigma)$
5. $\sup_{\theta \in N} \left\| \frac{\partial Q_n(\theta)^2}{\partial \theta \partial \theta'} - \frac{\partial Q(\theta)^2}{\partial \theta \partial \theta'} \right\| \xrightarrow{P} 0$
6. $\mathcal{H} \equiv \frac{\partial Q(\theta_0)^2}{\partial \theta \partial \theta'}$ is nonsingular

Theorem 6.6. Extremum Estimation Asymptotics Under Assumptions 6.1,

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, \mathcal{H}^{-1} \Sigma \mathcal{H}^{-1})$$

Proof. By the definition of $\hat{\theta}$ and the first two assumptions, we have that almost surely

$$\frac{\partial Q_n(\hat{\theta})}{\partial \theta} = 0$$

Using the third assumption and the Mean Value Theorem, we get

$$\frac{\partial Q_n(\hat{\theta})}{\partial \theta} = \frac{\partial Q_n(\theta_0)}{\partial \theta} + \frac{\partial Q_n(\bar{\theta})^2}{\partial \theta \partial \theta'} (\hat{\theta} - \theta_0)$$

where $\bar{\theta}$ is coordinate-wise between $\hat{\theta}$ and θ_0 , in particular $\bar{\theta} \xrightarrow{p} \theta_0$. From here, we can combine to rearrange and we want to show that:

$$\sqrt{n}(\hat{\theta} - \theta_0) = - \underbrace{\left(\frac{\partial Q_n(\bar{\theta})^2}{\partial \theta \partial \theta'} \right)^{-1}}_{\xrightarrow{p} \mathcal{H}^{-1}} \underbrace{\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta}}_{\xrightarrow{d} \mathcal{N}(0, \Sigma)}$$

The only thing left to show is the convergence to \mathcal{H} . To keep everything neat, define $H(\theta) = \frac{\partial Q(\theta)^2}{\partial \theta \partial \theta'}$ and $H_n(\cdot)$ analogously. Then we have that

$$\begin{aligned} \|H_n(\bar{\theta}) - \mathcal{H}\| &= \|H_n(\bar{\theta}) - H(\bar{\theta}) + H(\bar{\theta}) - \mathcal{H}\| \\ &\text{(by } \triangle) \leq \|H_n(\bar{\theta}) - H(\bar{\theta})\| + \|H(\bar{\theta}) - \mathcal{H}\| \\ &\text{(by Assumption 1)} \leq \sup_{\theta \in N} \|H_n(\theta) - H(\theta)\| + \|H(\bar{\theta}) - \mathcal{H}\| \\ &\text{(by Assumptions 3 \& 5)} \xrightarrow{p} 0 \end{aligned}$$

The claim now follows directly from nonsingularity of \mathcal{H} and Continuous Mapping Theorem. \square

Remark. We can slightly improve on this if the application is specifically nonlinear GMM. Recall that $\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \{\bar{g}_n(\theta)' W \bar{g}_n(\theta)\}$.

Assumption 6.2. We slightly refine the above assumptions to fit nonlinear GMM:

1. $\hat{\theta} \xrightarrow{p} \theta_0$
2. $\theta_0 \in \operatorname{int}(\Theta)$
3. $g(W; \theta)$ is almost surely continuously differentiable in an open neighborhood N of θ_0
4. $\sqrt{n} \bar{g}_n(\theta_0) \xrightarrow{d} \mathcal{N}(0, S)$ where S is positive definite
5. $\sup_{\theta \in N} \left\| \frac{\partial \bar{g}_n(\theta)}{\partial \theta'} - \mathbb{E} \left(\frac{\partial \bar{g}_n(\theta_0)}{\partial \theta'} \right) \right\| \xrightarrow{p} 0$
6. $\mathcal{G} \equiv \frac{\partial g(\theta_0)}{\partial \theta}$ is of full column rank

Theorem 6.7. Nonlinear GMM Asymptotics Under Assumptions 6.2,

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, (\mathcal{G}' W \mathcal{G})^{-1} \mathcal{G}' W S W \mathcal{G} (\mathcal{G}' W \mathcal{G})^{-1})$$

Remark. Two-stage efficient GMM works just as before! The main improvement is that we only need once-differentiability of $g(\cdot)$.

Example. Maximum Likelihood MLE is an extremely important special case. Say that we are able to specify

the distribution of data up to θ , so say the data are distributed with density

$$f(W_1, \dots, W_n; \theta)$$

where the function $f(\cdot)$ is known. Then the maximum likelihood estimator is

$$\hat{\theta}_{ML} = \operatorname{argmax}_{\theta \in \Theta} f(w_1, \dots, w_n; \theta)$$

Intuitively, this is the parameter value that maximizes the likelihood of observing the data that were actually observed. For discussion of maximum likelihood, we will think of extremum estimators as maximizing $Q(\cdot)$.

If we assume the data are i.i.d., we have the simplification

$$\begin{aligned} \hat{\theta}_{ML} &\equiv \operatorname{argmax}_{\theta \in \Theta} f(w_1, \dots, w_n; \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \prod_{i=1}^n f(w_i; \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log f(w_i; \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log f(w_i; \theta) \end{aligned}$$

This is a much easier (and often the only realistic) objective to compute. This is typically consistent even if the data are not i.i.d. (however, we will almost always assume that). The last step just serves to remind us that it is an m -estimator.

Definition. There are many different ways to think about *identification*. Here are three we have used:

- In linear moment-based models, it is a *rank condition*
- In extremum estimation, it is that θ_0 *uniquely minimizes* $Q(\cdot)$
- In Maximum Likelihood, it is

$$\theta \neq \theta_0 \implies \mathbb{P}\{f(W; \theta) \neq f(W; \theta_0)\} > 0$$

or equivalently,

$$\theta \neq \theta_0 \implies \exists A \in \mathcal{W}, \mathbb{P}\{A\} > 0, f(w; \theta) \neq f(w; \theta_0) \forall w \in A$$

where \mathcal{W} is the sample space (or the space of all possible realizations of W). Verbally, data that signal whether θ or θ_0 are true have positive probability.

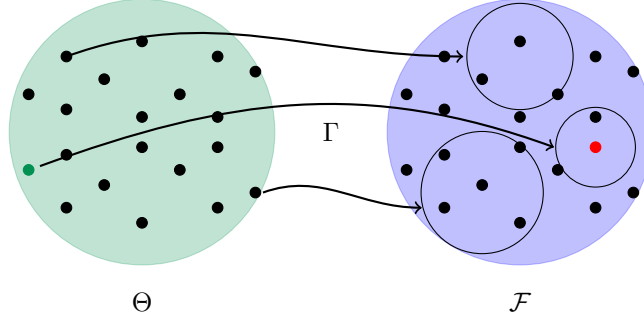
Question. What do these all have in common? **If we knew the population distribution of the data, we could back out θ_0 .**

- In linear moment-based models, the rank condition implies that the population moment conditions could be solved for θ_0
- In extremum estimation, uniqueness of the minimum at θ_0 means that knowledge of $Q(\cdot)$ implies knowledge of θ_0 (at least in principle)
- In Maximum Likelihood... we'll see.

Remark. The term identification is loaded. See [Lewbel \(2019\)](#) for all the (26!) different ways people use it. In empirical work, they often ask the question ‘where does your identification come from?’ This is unrelated to our usage, which corresponds also to the term *identifiability* in statistics.

Definition. We will introduce the following notation (motivated by [Matzkin's handbook](#) chapter on identification): (i) \mathcal{F} is the set of all possible population distributions of the data W ; (ii) Θ is the parameter space; and (iii) the correspondence $\Gamma : \Theta \rightarrow \mathcal{F}$ maps each parameter value to the set of distributions consistent with it.

We can think of the picture:



Remark. If likelihood is specified, $\Gamma(\cdot)$ is a singleton. In GMM, we have that $\Gamma(\theta) = \{F(W) \in \mathcal{F} : \mathbb{E}_F g(W; \theta) = 0\}$, so θ_0 is identified if $F \in \Gamma(\theta_0)$ implies that $\Gamma^{-1}(F) = \{\theta_0\}$. We usually consider θ identified if the above holds for all possible true values.

Remark. This can be used to motivate some extensions, not pursued here:

1. *Partial Identification:* $\Gamma^{-1}(F_0)$ is a set. This is completely uninformative if it is Θ , and point-identifying if it is $\{\theta_0\}$. Often it is somewhere in between.
2. *Irregular Identification or Ill-posed Inverse Problems:* $\Gamma^{-1}(\cdot)$ is sufficiently ill-behaved so that identifiability formally obtains but, for example, convergence of the empirical distribution F_n to F may imply convergence of $\Gamma^{-1}(F_n)$ to $\Gamma^{-1}(F)$ at a slower, if any, rate.

Remark. In many cases, the distribution of regressors X is not informative about θ . That is, we can write

$$f(Y, X; \theta) = f_y(Y | X; \theta) f_X(X)$$

then we have simplification

$$\begin{aligned} \hat{\theta}_{ML} &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log f(Y, X; \theta) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n (\log f_y(Y | X; \theta) + \log f_X(X)) \\ &= \operatorname{argmax}_{\theta \in \Theta} \sum_{i=1}^n \log f_y(Y | X; \theta) \end{aligned}$$

In practice, many ML estimators reflect this simplification. For the purpose of theoretical analysis, we always write the estimator as maximizing the complete likelihood.

Consistency of maximum likelihood follows from the m -estimator consistency result above. In fact:

Theorem 6.8. θ_0 uniquely maximizes $\mathbb{E} \log f(W; \theta)$ if and only if $\theta \neq \theta_0$ implies that $\mathbb{P}\{f(W; \theta) \neq f(W; \theta_0)\} > 0$.

Proof. We write

$$\begin{aligned}\mathbb{E} \log f(W; \theta) - \mathbb{E} \log f(W; \theta_0) &= \mathbb{E} \log \frac{f(W; \theta)}{f(W; \theta_0)} \leq \log \mathbb{E} \frac{f(W; \theta)}{f(W; \theta_0)} \\ &= \log \int \frac{f(w; \theta)}{f(w; \theta_0)} f(w; \theta_0) dw = \log \int f(w; \theta) dw = \log 1 = 0\end{aligned}$$

where the inequality is Jensen's Inequality, and it is strict unless

$$\frac{f(W; \theta)}{f(W; \theta_0)} \text{ constant almost surely} \iff \mathbb{P}\{f(W; \theta) \neq f(W; \theta_0)\} = 0$$

□

Remark. The structure of maximum likelihood allows us to both verify the “CLT assumption” and provide an expression for the asymptotic variance:

$$\begin{aligned}\int f(w; \theta) dw &= 1 \quad \forall \theta \in \Theta \\ \implies \int \frac{\partial f(w; \theta)}{\partial \theta} dw &= 0 \\ \implies \int \frac{\partial \log f(w; \theta)}{\partial \theta} f(w; \theta) dw &= 0 \\ \implies \mathbb{E} \left(\frac{\partial \log f(w; \theta)}{\partial \theta} \right) &= 0\end{aligned}$$

This result, called the *score equation*, is important in its own right: it tells us that maximum likelihood can be interpreted as a method of moments estimator. Taking derivatives again, we get:

$$\begin{aligned}\int \frac{\partial^2 \log f(w; \theta)}{\partial \theta \partial \theta'} f(w; \theta) dw + \int \frac{\partial \log f(w; \theta)}{\partial \theta} \frac{\partial \log f(w; \theta)}{\partial \theta'} f(w; \theta) dw &= 0 \\ \implies \mathbb{E} \left(\frac{\partial^2 \log f(w; \theta)}{\partial \theta \partial \theta'} \right) + \mathbb{E} \left(\frac{\partial \log f(w; \theta)}{\partial \theta} \frac{\partial \log f(w; \theta)}{\partial \theta'} \right) &= 0 \\ \implies \mathbb{E} \left(\frac{\partial^2 \log f(w; \theta)}{\partial \theta \partial \theta'} \right) &= -\mathbb{E} \left(\frac{\partial \log f(w; \theta)}{\partial \theta} \frac{\partial \log f(w; \theta)}{\partial \theta'} \right)\end{aligned}$$

The last line is, of course, the *information matrix equality*! Now we write

$$Q_n(\theta_0) = \frac{1}{n} \sum_{i=1}^n \log f(w_i; \theta_0) \implies \frac{\partial Q_n(\theta_0)}{\partial \theta} = \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(w_i; \theta_0)}{\partial \theta}$$

But we just showed that $\mathbb{E} \left(\frac{\partial \log f(w; \theta_0)}{\partial \theta} \right) = 0$. We thus have from the CLT:

$$\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \xrightarrow{d} \mathcal{N} \left(0, \mathbb{E} \left(\frac{\partial \log f(w; \theta_0)}{\partial \theta} \frac{\partial \log f(w; \theta_0)}{\partial \theta'} \right) \right)$$

This establishes part (4) of Assumptions 6.1. Substituting these findings into the theorem, we get

$$\begin{aligned}\sqrt{n}(\hat{\theta} - \theta_0) &\xrightarrow{d} \mathcal{N}\left(0, \underbrace{\left(\mathbb{E}\left(\frac{\partial^2 \log f(w; \theta_0)}{\partial \theta \partial \theta'}\right)\right)}_{\mathcal{H}}^{-1} \underbrace{\mathbb{E}\left(\frac{\partial \log f(w; \theta_0)}{\partial \theta} \frac{\partial \log f(w; \theta_0)}{\partial \theta'}\right)}_{S = -\mathcal{H}} \underbrace{\left(\mathbb{E}(\cdot)\right)}_{\mathcal{H}}^{-1}\right) \\ &= \mathcal{N}(0, -\mathcal{H}^{-1})\end{aligned}$$

where we use the Information Matrix Equality. Under our i.i.d. assumption, \mathcal{H} is the *(Fisher) information matrix* $\mathbb{I}(\theta_0)$, so ML asymptotically attains the *Cramer-Rao lower bound*. In fact, it is known (we will not show it) that ML is asymptotically efficient in the sense that it has the smallest asymptotic variance in a large class of regular estimators. This creates a strong case for using ML, as long as you are willing to specify a likelihood and can compute the ML estimator.

Remark. Whenever we have a complete likelihood, we can use maximum likelihood, but we could also use GMM – knowledge of the likelihood implies knowledge of the moment conditions, definitely the score equations but possibly others. So could GMM match (or possibly beat) the performance of ML? No! we have directly that

$$\begin{aligned}(G'S^{-1}G)^{-1} - \mathbb{I}(\theta_0)^{-1} &\text{ is positive semi-definite} \\ (G'S^{-1}G)^{-1} &= \mathbb{I}(\theta_0)^{-1} \text{ if } g(w, \theta) = \frac{\partial \log f(w; \theta_0)}{\partial \theta}\end{aligned}$$

Thus, GMM cannot asymptotically beat ML estimation, as $\text{aVar}(\hat{\theta}_{GMM}) \geq \text{aVar}(\hat{\theta}_{ML})$. If the likelihood is known, GMM can trivially match ML by mimicking it, but since those moment conditions would reflect likelihood information, we cannot in general get ML efficiency without knowing the likelihood.

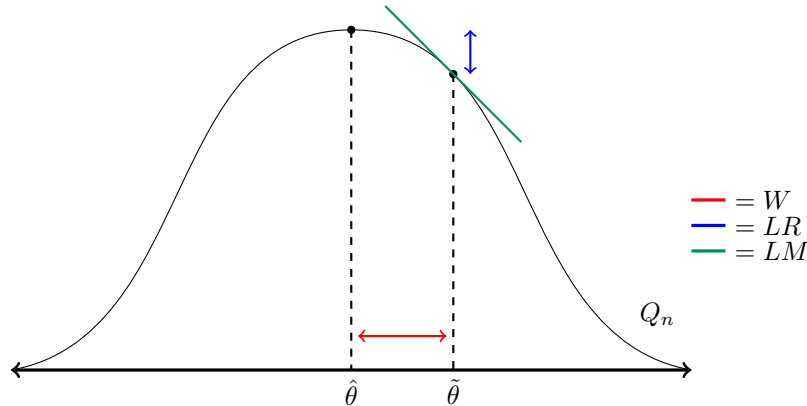
Example. Suppose we want to test $\mathbb{H}_0 : r(\theta) = 0$ where $r(\cdot)$ is a known function whose Jacobian $R(\cdot)$ is both continuous and has full rank at θ_0 . The *trinity* of test statistics are:

1. *Wald*: $W = nr(\hat{\theta})'(R(\hat{\theta})\hat{\Sigma}^{-1}R(\hat{\theta})')^{-1}r(\hat{\theta})$
2. *Likelihood Ratio*: $LR = 2n(Q_n(\hat{\theta}) - Q_n(\tilde{\theta}))$
3. *Lagrange Multiplier*: $LM = n\frac{\partial Q_n(\tilde{\theta})'}{\partial \theta}\tilde{\Sigma}^{-1}\frac{\partial Q_n(\tilde{\theta})}{\partial \theta}$

where $\tilde{\theta}$ is the *constrained estimator*

$$\tilde{\theta} \equiv \underset{\theta \in \Theta}{\operatorname{argmin}} Q_n(\theta) \text{ s.t. } r(\theta) = 0$$

and where $(\hat{\Sigma}, \tilde{\Sigma})$ estimate the outer product of gradients at $(\hat{\theta}, \tilde{\theta})$. We can illustrate the hypotheses as:



Assumption 6.3. We assume that:

1. $\sqrt{n}(\hat{\theta} - \theta_0) = -\mathcal{H}^{-1}\sqrt{n}\frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1)$
2. $\frac{\partial Q_n(\theta_0)}{\partial \theta} \xrightarrow{d} \mathcal{N}(0, \Sigma)$, with Σ positive definite
3. $\sqrt{n}(\tilde{\theta} - \theta_0) = O_P(1)$
4. $\Sigma = -\mathcal{H}$

Theorem 6.9. Under Assumptions 6.3, all of the tests (Wald, Likelihood Ratio, and Lagrange Multiplier) converge in distribution to $\chi^2_{\#r}$. Furthermore, they are asymptotically equivalent.

Remark. We will take on faith that $\sqrt{n}(\tilde{\theta} - \theta_0) = O_P(1)$, Hayashi shows it from the theory of constrained estimators. We will only spell out the details for maximum likelihood, for everything else \mathcal{H} must be redefined. We do need that $\mathcal{H} = -\Sigma$, meaning that ML is well-specified. We will come back to misspecified models.

Proof. (Only of the convergence statement). The argument for the Wald statistic is entirely above. The first order condition for the constrained estimation problem can be written as

$$\sqrt{n}\frac{\partial Q_n(\tilde{\theta})}{\partial \theta} + \sqrt{n}R(\tilde{\theta})'\gamma_n = 0 \quad ; \quad \sqrt{n}r(\tilde{\theta}) = 0$$

We can use the Mean Value Theorem to write

$$\begin{aligned} r(\tilde{\theta}) &= r(\theta_0) + R(\tilde{\theta})(\tilde{\theta} - \theta_0) \\ \implies \sqrt{n}r(\tilde{\theta}) &= \sqrt{n}R(\tilde{\theta})(\tilde{\theta} - \theta_0) \\ &= \underbrace{\sqrt{n}(R(\tilde{\theta}) - R(\theta_0))(\tilde{\theta} - \theta_0)}_{\xrightarrow{p} 0} + \sqrt{n}R(\theta_0)(\tilde{\theta} - \theta_0) \\ &= R(\theta_0) \cdot \sqrt{n}(\tilde{\theta} - \theta_0) + o_P(1) \end{aligned}$$

Next, a Taylor expansion of $\frac{\partial Q_n(\theta)}{\partial \theta}$ about θ_0 yields

$$\sqrt{n}\frac{\partial Q_n(\tilde{\theta})}{\partial \theta} = \underbrace{\sqrt{n}\frac{\partial Q_n(\theta_0)}{\partial \theta}}_{\xrightarrow{d} \mathcal{N}(0, \Sigma)} + \underbrace{\sqrt{n}\frac{\partial^2 Q_n(\theta_0)}{\partial \theta \partial \theta'}}_{\xrightarrow{p} \mathcal{H}}(\tilde{\theta} - \theta_0) + o_P(1)$$

The second and third assumptions now imply that $\sqrt{n}\frac{\partial Q_n(\tilde{\theta})}{\partial \theta}$, and thus $\sqrt{n}\gamma_n$, are of order $O_P(1)$. This in turn allows us to write

$$R(\tilde{\theta})'\sqrt{n}\gamma_n = R(\theta_0)'\sqrt{n}\gamma_n + (R(\tilde{\theta}) - R(\theta_0))'\sqrt{n}\gamma_n = R(\theta_0)'\sqrt{n}\gamma_n + o_P(1)$$

Next, collecting of terms. We have:

$$\begin{aligned} \sqrt{n}r(\tilde{\theta}) &= 0 \\ \sqrt{n}r(\tilde{\theta}) &= R(\theta_0) \cdot \sqrt{n}(\tilde{\theta} - \theta_0) + o_P(1) \\ \implies R(\theta_0) \cdot \sqrt{n}(\tilde{\theta} - \theta_0) &= o_P(1) \end{aligned}$$

as well as:

$$\begin{aligned}
& \sqrt{n} \frac{\partial Q_n(\tilde{\theta})}{\partial \theta} + \sqrt{n} R(\tilde{\theta})' \gamma_n = 0 \\
& \sqrt{n} \frac{\partial Q_n(\tilde{\theta})}{\partial \theta} = \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + \mathcal{H} \sqrt{n}(\tilde{\theta} - \theta_0) + o_P(1) \\
& R(\tilde{\theta})' \sqrt{n} \gamma_n = R(\theta_0)' \sqrt{n} \gamma_n + o_P(1) \\
& \implies \mathcal{H} \sqrt{n}(\tilde{\theta} - \theta_0) + R(\theta_0)' \sqrt{n} \gamma_n = -\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1)
\end{aligned}$$

This implicitly characterizes the joint distribution, but we can make it explicit, by consolidating into

$$\begin{bmatrix} \mathcal{H} & R' \\ R & 0 \end{bmatrix} \cdot \begin{bmatrix} \sqrt{n}(\tilde{\theta} - \theta_0) \\ \sqrt{n} \gamma_n \end{bmatrix} = \begin{bmatrix} -\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \\ 0 \end{bmatrix} + o_P(1)$$

which implies that

$$\sqrt{n} \begin{bmatrix} \tilde{\theta} - \theta_0 \\ \gamma_n \end{bmatrix} = \begin{bmatrix} -\mathcal{H}^{-1} + \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \\ -(R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \end{bmatrix} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1)$$

This gets us the LM statistic fairly quickly:

$$\begin{aligned}
\sqrt{n} \gamma_n &= -(R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1) \\
&\xrightarrow{d} \mathcal{N}(0, (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \Sigma \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1}) \\
&= \mathcal{N}(0, (R \Sigma^{-1} R')^{-1}) \\
&\implies \sqrt{n} \gamma_n' R \Sigma^{-1} R' \sqrt{n} \gamma_n \xrightarrow{d} \chi_{\#r}^2
\end{aligned}$$

We conclude with another MVT expansion:

$$Q_n(\tilde{\theta}) = Q_n(\hat{\theta}) + \frac{\partial Q_n(\hat{\theta})}{\partial \theta} (\tilde{\theta} - \hat{\theta}) + \frac{1}{2} (\tilde{\theta} - \hat{\theta})' \frac{\partial^2 Q_n(\bar{\theta})}{\partial \theta \partial \theta'} (\tilde{\theta} - \hat{\theta})$$

However, the first partial of $Q_n(\cdot)$ is zero (almost surely) and the second partial converges to the (negative) curvature. Thus,

$$\begin{aligned}
2n(Q_n(\hat{\theta}) - Q_n(\tilde{\theta})) &= -\sqrt{n}(\tilde{\theta} - \hat{\theta})(\mathcal{H} + o_P(1))\sqrt{n}(\tilde{\theta} - \hat{\theta}) \\
&= -\sqrt{n}(\tilde{\theta} - \hat{\theta})\mathcal{H}\sqrt{n}(\tilde{\theta} - \hat{\theta}) + o_P(1)
\end{aligned}$$

We can substitute from the matrix equation to get

$$\begin{aligned}
\sqrt{n}(\tilde{\theta} - \hat{\theta}) &= -(\mathcal{H}^{-1} - \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1}) \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + \mathcal{H}^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1) \\
&= \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} + o_P(1)
\end{aligned}$$

This determines the distribution of the LR statistic. The rest is algebra:

$$\begin{aligned}
2n(Q_n(\hat{\theta}) - Q_n(\tilde{\theta})) &\approx -\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} (\mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1})' \mathcal{H} \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \\
&= -\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \mathcal{H}^{-1} R' (R \mathcal{H}^{-1} R')^{-1} R \mathcal{H}^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \\
&= \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \Sigma^{-1} R' (R \Sigma^{-1} R')^{-1} R \Sigma^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta}
\end{aligned}$$

Recalling that $\sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} \xrightarrow{d} \mathcal{N}(0, \Sigma)$, we have

$$\begin{aligned}
R \Sigma^{-1} \sqrt{n} \frac{\partial Q_n(\theta_0)}{\partial \theta} &\xrightarrow{d} \mathcal{N}(0, R \Sigma^{-1} \Sigma \Sigma^{-1} R') = \mathcal{N}(0, R \Sigma^{-1} R') \\
&\implies LR \xrightarrow{d} \chi_{\#r}^2
\end{aligned}$$

□

Remark. Why is it called the likelihood ratio statistic? If we take the likelihood literally, then

$$n(Q_n(\hat{\theta}) - Q_n(\tilde{\theta})) = \sum_{i=1}^n \ell(\hat{\theta}) - \sum_{i=1}^n \ell(\tilde{\theta}) = \frac{f(W_1, \dots, W_n; \hat{\theta})}{f(W_1, \dots, W_n; \tilde{\theta})}$$

The additional factor of 2 aligns the statistic with the others, but the interpretation of $Q_n(\cdot)$ is not essential.

7 Worked Examples

Example. *Poisson* Let Y be distributed i.i.d. Poisson with true parameter $\lambda_0 > 0$. Recall that $\mathbb{P}\{Y = k\} = \lambda^k e^{-\lambda} / k!$, with mean and variance λ . Intuitively, we of course will estimate λ using \bar{y} , and the Lindeberg-Lévy CLT (see Econometrics I notes for a full treatment) immediately yields $\sqrt{n}(\bar{y} - \lambda_0) \xrightarrow{d} \mathcal{N}(0, \lambda_0)$. To analyze this as an example of maximum likelihood, we write:

$$\begin{aligned}
Q_n(\lambda) &= \frac{1}{n} \sum_i (y_i \log \lambda - \lambda - \log(y_i!)) \\
\frac{\partial Q_n(\lambda)}{\partial \lambda} &= \frac{1}{n} \sum_i \left(\frac{y_i}{\lambda} - 1 \right) \\
\frac{\partial^2 Q_n(\lambda)}{\partial \lambda^2} &= -\frac{1}{n} \sum_i \frac{y_i}{\lambda^2} \\
Q(\lambda) &= \mathbb{E}[Y \log \lambda - \lambda - \log Y!] \\
\frac{\partial Q(\lambda)}{\partial \lambda} &= \mathbb{E} \left[\frac{Y}{\lambda} - 1 \right] \\
\frac{\partial^2 Q(\lambda)}{\partial \lambda^2} &= -\mathbb{E} \left[\frac{Y}{\lambda^2} \right]
\end{aligned}$$

We can see that we have strict concavity, so the maximum likelihood estimator is characterized by the first order condition. In particular,

$$\frac{1}{n} \sum_i \left(\frac{y_i}{\lambda} - 1 \right) = 0 \implies \hat{\lambda}_{ML} = \bar{y}$$

We can next apply the consistency theorem for extremum estimators. Pointwise convergence of Q_n to Q is clear. The parameter space is not compact but since Q_n is strictly concave and λ is a scalar, that is not necessary. Of the conditions for asymptotic normality, the first was established above, the second holds as long as $\lambda > 0$, the third is obvious, the fourth holds because Y is i.i.d. with first and second moments λ_0 , meaning that

$$\sqrt{n} \frac{\partial Q_n(\lambda_0)}{\partial \lambda} = \sqrt{n} \frac{1}{n} \sum_i \left(\frac{Y_i}{\lambda_0} - 1 \right) \xrightarrow{d} \mathcal{N}(0, \lambda_0^{-1})$$

The fifth condition follows directly, the sixth we have directly because $\frac{\partial^2 Q(\lambda_0)}{\partial \lambda^2} = -\frac{1}{\lambda_0}$, and the last because $\lambda_0 > 0$. Thus, the consistency theorem applies, and we have that

$$\sqrt{n}(\hat{\lambda} - \lambda_0) \xrightarrow{d} \mathcal{N}(0, (-\lambda_0^{-1})^{-1} \lambda_0^{-1} (-\lambda_0^{-1})^{-1}) = \mathcal{N}(0, \lambda_0)$$

Note that these results do not apply if $\lambda_0 = 0$, because interiority is violated. In this case, we have an actual failure of result – the distribution of both Y and \bar{y} will be degenerate at 0. The above result does not uniformly hold as $\lambda_0 \rightarrow 0$; it fails along drifting parameters of the form $\lambda_n = \gamma/\sqrt{n}$.

Example. *Binary Response* Binary response models can generally be expressed in the form

$$Y = \mathbb{1}\{\phi(X, \varepsilon; \theta_0) \geq 0\}$$

which is frequently specialized to

$$Y = \mathbb{1}\{X'\beta - \varepsilon \geq 0\} \iff \mathbb{P}\{Y = 1 \mid X\} = \mathbb{P}\{\varepsilon \leq X'\beta\} = F_\varepsilon(X'\beta)$$

Assume that $\mathbb{E}XX'$ is nonsingular and that F_ε is strictly increasing. Then the model is identified up to scale normalization and (if X has a constant) a location normalization. Other than that, different assumptions about F_ε lead to different models. If we assume that ε is logistically distributed, we have the logit model

$$\mathbb{P}\{Y = 1 \mid X\} = \frac{\exp(X'\beta)}{1 + \exp(X'\beta)}$$

Which admits log likelihood

$$\log f(Y \mid X; \theta) = Y \log \frac{\exp(X'\beta)}{1 + \exp(X'\beta)} + (1 - Y) \log \frac{1}{1 + \exp(X'\beta)}$$

and the maximum likelihood estimator is characterized as a maximizer of

$$\begin{aligned} Q_n(\beta) &= \frac{1}{n} \sum_i \left(Y_i \log \frac{\exp(X'_i \beta)}{1 + \exp(X'_i \beta)} + (1 - Y_i) \log \frac{1}{1 + \exp(X'_i \beta)} \right) \\ &= \frac{1}{n} \sum_i (Y_i X'_i \beta - Y_i \log(1 + \exp(X'_i \beta)) - (1 - Y_i) \log(1 + \exp(X'_i \beta))) \\ &= \frac{1}{n} \sum_i (Y_i X'_i \beta - \log(1 + \exp(X'_i \beta))) \end{aligned}$$

We can thus write

$$\begin{aligned} \frac{\partial Q_n(\beta)}{\partial \beta} &= \frac{1}{n} \sum_i \left(Y_i X_i - \frac{\exp(X'_i \beta)}{1 + \exp(X'_i \beta)} X_i \right) = \frac{1}{n} \sum_i (Y_i - F_\varepsilon(X'_i \beta)) X_i \\ \frac{\partial^2 Q_n(\beta)}{\partial \beta \partial \beta'} &= -\frac{1}{n} \sum_i F_\varepsilon(X'_i \beta) (1 - F_\varepsilon(X'_i \beta)) X_i X'_i \end{aligned}$$

where we apply the formula $F(t) = e^t/(1 + e^t) \implies F'(t) = F(t)(1 - F(t))$ to the CDF. As in the first example, we could directly attain consistency and asymptotic normality by checking the conditions, which all hold.

Remark. If we replaced sample averages with expectations in the above, we would have

$$\begin{aligned}\frac{\partial Q(\beta)}{\partial \beta} &= \mathbb{E}((Y - F_\varepsilon(X'\beta))X) \\ \frac{\partial^2 Q(\beta)}{\partial \beta \partial \beta'} &= -\mathbb{E}(F_\varepsilon(X'\beta)(1 - F_\varepsilon(X'\beta))XX')\end{aligned}$$

At first glance, the information matrix equality may not appear obvious. However, we can convert this to:

$$\begin{aligned}& \mathbb{E} \left[((Y - F_\varepsilon(X'\beta))X) ((Y - F_\varepsilon(X'\beta))X)' \right] \\ &= \mathbb{E} \left[(Y - F_\varepsilon(X'\beta))^2 XX' \right] \\ &= \mathbb{E} \left[\mathbb{E}[(Y - F_\varepsilon(X'\beta))^2 | X] XX' \right] \\ &= \mathbb{E}[\text{Var}(Y | X) XX'] \\ &= \mathbb{E}[F_\varepsilon(X'\beta)(1 - F_\varepsilon(X'\beta)) XX'] \\ &= -\frac{\partial^2 Q(\beta)}{\partial \beta \partial \beta'}\end{aligned}$$

where we use the Law of Iterated Expectations and the knowledge that Y is distributed Bernoulli once we condition on X .

Remark. We could also estimate this model using GMM. To do so, we need a function g of the dimensionality of X such that $\mathbb{E}g(Y, X; \beta) = 0$. As a general rule, if conditional expectations can be written out, they immediately give rise to such functions. In the current example, from LIE we get

$$\mathbb{E}(X(Y - F_\varepsilon(X'\beta))) = \mathbb{E}(X\mathbb{E}(Y - F_\varepsilon(X'\beta) | X)) = 0$$

so our GMM estimator is defined by the sample analog

$$\frac{1}{n} \sum_i X_i(Y_i - F_\varepsilon(X_i'\hat{\beta})) = 0$$

So the estimators algebraically coincide! The natural GMM estimator uses the score equations as moment conditions and therefore is exactly the ML estimator – replicating even the variance precisely.

Remark. If we consider $\mathbb{E}[Y|X] = G(X'\beta)$, we consider $G(\cdot)$ the *link function* and $X'\beta$ the *linear index*. If we shift this by a sigmoid, we can convert this to logit or probit models quite easily. For the probit, we have that

$$\mathbb{E}[Y | X] = \Phi(X'\beta) \iff Y = \mathbb{1}\{X'\beta - \varepsilon \geq 0\}, \varepsilon \sim \mathcal{N}(0, 1)$$

In logit and probit models, coefficients are in general sign interpretable; relative absolute values are *sometimes* interpretable; absolute values without context are not interpretable. We have that

$$\frac{\partial \mathbb{E}(Y | X)}{\partial X_j} = \beta_j \cdot g(X'\beta)$$

where $g(\cdot) \equiv G'(\cdot)$. The sign of this term will be the sign of the coefficient, as $g(\cdot)$ is positive always, but the marginal effect depends on the magnitude of $g(\cdot)$ which is not in general known. Put plainly, the marginal effect of X_j on Y is now described by two numbers, and depends on where X_j is. If you *need* an interpretable number, there are some options: marginal effect at a particular fixed X , average estimated effect which is $\mathbb{E}_n \hat{\beta}_j \cdot g(X'\hat{\beta})$, and the estimated effect at the average which is $\hat{\beta}_j \cdot g(\mathbb{E}_n X'\hat{\beta})$.

Example. *Type II Tobit* The following model is known as *Type II Tobit*:

$$\begin{aligned} Y_1^* &= X_1' \beta_1 + \varepsilon_1 \\ Y_2^* &= X_2' \beta_2 + \varepsilon_2 \\ Y_1 &= \mathbb{1}\{Y_1^* \geq 0\} \\ Y_2 &= Y_2^* \cdot \mathbb{1}\{Y_1^* \geq 0\} \\ \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix} &\sim \mathcal{N}\left(0, \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix}\right) \end{aligned}$$

Note that you observe only Y_1 and Y_2 , not the respective latent starred variables. The motivation here is that the equation of interest is for Y_2^* , called the *outcome equation*. The equation for Y_1^* is called the *selection equation*. The parameter of interest is β_2 . If we observed Y_2^* , we could estimate it by straight OLS. However, since we are selecting for when $Y_1^* \geq 0$, we are selecting for high values of ε_1 . As long as $\rho \neq 0$, the correlation between ε_1 and ε_2 means that we have selection also in the sample of interest.

The motivating examples here are typically that the decision to enter some market is dependent on a threshold that does not affect the values in the market – think if the decision to enter the labor force is dependent on having children, but wages are not directly determined by children, or if foreign aid is dependent on a certain political stance but the level of foreign aid does not depend on said political stance.

The selection model in the first and third equations gives us the likelihood

$$\mathbb{P}\{Y_1 = 0 \mid X_1, X_2; \theta\} = \Phi\left(-\frac{X_1' \beta_1}{\sigma_1}\right)$$

Thus, the likelihood of Y_2 conditional on X_1, X_2 , and $Y_1 = 1$ (meaning Y_2 is observed at all) is:

$$\begin{aligned} f(Y_2 \mid X_1, X_2, Y_1 = 1; \theta) &= \mathbb{P}\{Y_1 = 1 \mid X_1, X_2; \theta\} \cdot f(Y_2 \mid X_1, X_2; \theta) / \Phi\left(\frac{X_1' \beta_1}{\sigma_1}\right) \\ &= \frac{1}{\sigma_2} \Phi\left(\frac{X_1' \beta_1 + \rho \frac{\sigma_1}{\sigma_2} (Y_2 - X_2' \beta_2)}{\sigma_1 \sqrt{1 - \rho^2}}\right) \phi\left(\frac{Y_2 - X_2' \beta_2}{\sigma_2}\right) / \Phi\left(\frac{X_1' \beta_1}{\sigma_1}\right) \end{aligned}$$

where we used the (conditional) distributions of the errors.

We will establish identification using the thought experiment: If we had perfect knowledge of the above likelihoods, could we back out the true parameter values? Actually, the answer is no. The probit part of the model identifies the ratio β_1/σ_1 , but neither parameter on their own. The natural way to get identification is to normalize $\sigma_1 = 1$, and work to identify all other parameters. This gives us the simplification:

$$\begin{aligned} \mathbb{P}\{Y_1 = 0 \mid X_1, X_2; \theta\} &= \Phi(-X_1' \beta_1) \\ f(Y_2 \mid X_1, X_2, Y_1 = 1; \theta) &= \frac{1}{\sigma_2} \Phi\left(\frac{X_1' \beta_1 + \frac{\rho}{\sigma_2} (Y_2 - X_2' \beta_2)}{\sqrt{1 - \rho^2}}\right) \phi\left(\frac{Y_2 - X_2' \beta_2}{\sigma_2}\right) / \Phi(X_1' \beta_1) \end{aligned}$$

With $\sigma_1 = 1$, our selection equation now identifies β_1 , so we will treat $\Phi(X_1' \beta_1)$ as known. Define:

$$\tilde{f}(Y_2 \mid X_1, X_2; \theta) = \frac{1}{\sigma_2} \Phi\left(\frac{X_1' \beta_1 + \frac{\rho}{\sigma_2} (Y_2 - X_2' \beta_2)}{\sqrt{1 - \rho^2}}\right) \phi\left(\frac{Y_2 - X_2' \beta_2}{\sigma_2}\right)$$

The partials are:

$$\begin{aligned}\frac{\partial \tilde{f}(\cdot)}{\partial Y_2} &= \frac{1}{\sigma_2^2} \Phi(\cdot) \phi'(\cdot) + \frac{\rho}{\sigma_2^2 \sqrt{1-\rho^2}} \Phi'(\cdot) \phi(\cdot) \\ \frac{\partial \tilde{f}(\cdot)}{\partial X_2} &= -\frac{\beta_2}{\sigma_2^2} \Phi(\cdot) \phi'(\cdot) - \frac{\beta_2 \rho}{\sigma_2^2 \sqrt{1-\rho^2}} \Phi'(\cdot) \phi(\cdot) = -\beta_2 \cdot \frac{\partial \tilde{f}(\cdot)}{\partial Y_2} \\ \frac{\partial \tilde{f}(\cdot)}{\partial X_1} &= \frac{\beta_1}{\sigma_2 \sqrt{1-\rho^2}} \Phi'(\cdot) \phi(\cdot)\end{aligned}$$

First, note that β_2 can be identified through the partials with respect to X_2 and Y_2 . Having identified β_1 and β_2 , we can choose to evaluate \tilde{f} at arguments where $X_1' \beta_1 = Y_2 - X_2' \beta_2 = 0$. At these arguments, we will have

$$\begin{aligned}\frac{\partial \tilde{f}(\cdot)}{\partial X_1} &= \frac{\beta_1}{\sigma_2 \sqrt{1-\rho^2}} (\phi(0))^2 \\ \frac{\partial \tilde{f}(\cdot)}{\partial Y_2} &= \frac{\rho}{\sigma_2^2 \sqrt{1-\rho^2}} (\phi(0))^2\end{aligned}$$

Since β_1 is known, this is essentially two equations in the two remaining unknowns ρ and σ_2 . We finish with two remarks:

Remark. This is an example of a non-constructive identification proof: we established that perfect knowledge of the likelihood would allow us to back out the parameter values, but our identification strategy should not be to solve sample analogs of these equations, since they involve evaluating estimated derivatives of estimated densities at estimated parameter values – we have a *ton* of noise.

Remark. This argument made use of some support conditions, where in the thought experiment we freely took derivatives of likelihoods and used them. However, we can only evaluate these derivatives at values of (X_1, X_2, Y_1, Y_2) on the support of the true distribution. Because of the normality assumption on ε_2 , we know that $Y_2 - X_2' \beta_2 = 0$ occurs on the support. We do not actually know the same for $X_1' \beta_1 = 0$, but that assumption could be relaxed as long as X_1 has *some* variation.

So how should we actually estimate this model? The obvious approach is maximum likelihood, where we have the objective:

$$Q_n(\theta) = \frac{1}{n} \sum_{i=1}^n \left[(1 - Y_i) \log \Phi(-X_i' \beta_1) + Y_i \left(\log \Phi \left(\frac{X_i' \beta_1 + \frac{\rho}{\sigma_2} (Y_i - X_i' \beta_2)}{\sqrt{1-\rho^2}} \right) - \frac{1}{2} (Y_i - X_i' \beta_2)^2 - \log \sigma_2 \right) \right]$$

where $\theta = (\beta_1, \beta_2, \rho, \sigma_2)$.

Example. “*Heckit*” ML is the statistically efficient way to estimate these models, and a lot of implementations exist, with some shortcuts. However, as the dimensions of the parameter vectors increase, this is an involved problem because the likelihood is multimodal. Heckman proposed a two-step method, that follows the observation that

$$\mathbb{E}(Y_2 \mid X_1, X_2, Y_1 = 1) = \mathbb{E}(X_2' \beta_2 + \varepsilon_2 \mid X_1, X_2, Y_1 = 1) = X_2' \beta_2 + \mathbb{E}(\varepsilon_2 \mid \varepsilon_1 \geq -X_1' \beta_1)$$

Recall that, if we assume z_i is standard normal, the following holds:

$$\begin{aligned}\mathbb{E}(z_i \mid z_i \geq t) &= \frac{\int_{z=t}^{\infty} z \phi(z) \frac{\partial z}{\partial z}}{\int_{z=t}^{\infty} \phi(z) \frac{\partial z}{\partial z}} = \frac{(2\pi)^{-1/2} \int_{z=t}^{\infty} z \exp(-z^2/2) \frac{\partial z}{\partial z}}{\Phi(-t)} \\ &= \frac{(2\pi)^{-1/2} [-\exp(-z^2/2)]_t^{\infty}}{\Phi(-t)} = \frac{(2\pi)^{-1/2} \exp(-t^2/2)}{\Phi(-t)} = \frac{\phi(t)}{\Phi(-t)} := \lambda(-t)\end{aligned}$$

where the last equality defines the *Inverse Mills Ratio*. Thus, in our terms we have

$$\begin{aligned}\mathbb{E}(\varepsilon_2 \mid \varepsilon_1 \geq -X_1'\beta_1) &= \frac{\int_{-X_1'\beta_1}^{\infty} \mathbb{E}(\varepsilon_2 \mid \varepsilon_1)\phi(\varepsilon_1) \partial\varepsilon_1}{\int_{-X_1'\beta_1}^{\infty} \phi(\varepsilon_1) \partial\varepsilon_1} \\ &= \frac{\int_{-X_1'\beta_1}^{\infty} \rho\sigma_2\varepsilon_1\phi(\varepsilon_1) \partial\varepsilon_1}{\int_{-X_1'\beta_1}^{\infty} \phi(\varepsilon_1) \partial\varepsilon_1} = \rho\sigma_2\lambda(X_1'\beta_1)\end{aligned}$$

Thus, when $Y_1 = 1$, we could write

$$Y_2 = X_2'\beta_2 + \rho\sigma_2\lambda(X_1'\beta_1) + \eta_i$$

where $\mathbb{E}(\eta_i \mid X_1, X_2, Y_1 = 1) = 0$. If we knew β_1 , we could estimate β_2 using OLS of Y_2 on $(X_2, \lambda(X_1'\beta_1))$. In reality, things are a bit more complicated because β_1 must also be estimated. Heckman (1976) established that the following two-step procedure works:

Step 1. Use probit to estimate β_1 , call this estimator $\hat{\beta}_1$.

Step 2. Restrict attention to observations where $Y_1 = 1$. Use OLS to estimate the equation

$$Y_2 = X_2'\beta_2 + \rho\sigma_2\lambda(X_1'\hat{\beta}_1) + \eta$$

That this works is not obvious due to the estimated regressor on the right hand side, but it is nonetheless true. Unsurprisingly, OLS standard errors would not be valid.

The Heckit method easily generalizes to variations of the above model, and the model is coded into most statistical packages. Notice, however, that it is inefficient: it assumes normality for identification, and is sensitive to the normality assumption failure, but doesn't use them for estimation! In particular, an ML estimator would use second-stage information also in the estimation of β_1 . The choice between ML and Heckit depends on how complicated the likelihood is in a certain application. Both are implemented in Stata.

Example. (*Smoothed*) *Maximum Score* (from Manski (1975), Kim & Pollard (1990), Chamberlain (1986), and Horowitz (1992)) Consider the binary choice model

$$Y = \mathbb{1}\{X'\beta + \varepsilon \geq 0\}$$

where the researcher observes (Y, X) . We do not assume an exact distribution for ε , but we do assume that ε is continuous and $\text{med}(\varepsilon) = 0$. For simplicity, we also assume continuous X with full support rather than having a constant.

These assumptions are weaker than probit, so we will identify β only up to scale. Unlike earlier, it is convenient (and standard) to normalize $\|\beta\| = 1$, and impose the same restrictions on estimators. Then we have

$$\hat{\beta} \in \underset{b: \|b\|=1}{\operatorname{argmax}} \mathbb{E}((2Y - 1)\mathbb{1}\{X'\beta \geq 0\})$$

because by LIE,

$$\mathbb{E}((2Y - 1)\mathbb{1}\{X'\beta \geq 0\}) = \mathbb{E}[\mathbb{E}(2Y - 1 \mid X) \mathbb{1}\{X'\beta \geq 0\}]$$

and the right hand side outer integrand is maximized pointwise by β because

$$\mathbb{E}(2Y - 1 \mid X) \geq 0 \iff X'\beta \geq 0$$

If X has full support, the above argmax is unique. Without this assumption, we may have partial identification of β , which can be thought of as being able to ‘wiggle’ the separating hyperplane.

Note that the empirical distribution of X_i has at most n mass points, so it never has full support. Thus, the estimator

$$\hat{\beta} = \operatorname{argmax}_{b: \|b\|=1} \frac{1}{n} \sum_{i=1}^n (2Y_i - 1) \mathbb{1}\{X_i' \beta \geq 0\}$$

is not well-defined since the argmax is not unique. In the following, let $\hat{\beta}$ be a measurable selection from the argmax (e.g. the element that minimizes the first component, then the second, then the third, etc).

This is the *maximum score* estimator. It has historical importance as one of the first nonparametric¹³ estimators. It is also supremely ill-behaved. Under reasonable conditions, consistency can be established like the examples earlier. However, the asymptotics are horrible. The Hessian at the sample argmax is zero, which indicates something that is true – it converges at $n^{1/3}$, slower than \sqrt{n} . Specifically, an estimator with \sqrt{n} -consistency does not exist under the assumptions, and the asymptotic distribution is intractable.

These issues are basically entirely due to the non-smoothness of the objective function. Would smoothing the objective function get you a better-behaved estimator? Yes! Specifically, write:

$$\hat{\beta} = \operatorname{argmax}_{b: \|b\|=1} \frac{1}{n} \sum_{i=1}^n (2Y_i - 1) g_n(X_i' \beta)$$

where the function $g_n(\cdot)$ is smooth, has $g_n(0) = \frac{1}{n}$, and goes to 0 as the argument goes to $-\infty$, and 1 as the argument goes to ∞ . Specifically, we need that $g_n(t) \rightarrow \mathbb{1}\{t \geq 0\}$ at a certain rate as $n \rightarrow \infty$. This is very similar to kernel density estimation.

This estimator is called the *smoothed maximum score* estimator. It is asymptotically normal, and converges at a rate arbitrarily close to \sqrt{n} (assuming X_i has some reasonable properties and $g_n(\cdot)$ is chosen smartly).

Example. *Maximum of a Uniform Distribution* This is an example on which ML and GMM might disagree. Let $X \sim U[0, \alpha_0]$, where we seek to estimate α_0 . Note that $\mathbb{E}X = \alpha_0/2$ and so $\hat{\alpha} = 2\bar{X}$ is a GMM estimator based on the moment condition $\mathbb{E}(2X - \alpha_0) = 0$. We understand the behavior of this estimator extremely well. However, it is not in this example the ML estimator and ends up being quite inefficient.

We can compute the ML estimator. The likelihood for a single observation is $f(x; \alpha) = \frac{1}{\alpha} \mathbb{1}\{0 \leq x \leq \alpha\}$. The sample criterion function is

$$Q_n(\alpha) = \frac{1}{n} \sum_{i=1}^n \log \left(\frac{1}{\alpha} \mathbb{1}\{0 \leq x \leq \alpha\} \right) = \begin{cases} -\log \alpha & \max_i x_i \in [0, \alpha] \\ -\infty & \text{otherwise} \end{cases}$$

with population analog

$$Q(\alpha) = \mathbb{E} \log f(X; \alpha) = \begin{cases} -\log \alpha & \alpha \leq \alpha_0 \\ -\infty & \text{otherwise} \end{cases}$$

By inspection, these problems are solved by $\hat{\alpha}_0 = \max_i X_i$.

These objective functions are discontinuous and not differentiable, exactly at the optimum! None of our extremum estimator theorems are available, and it's a red flag for the asymptotics. In fact, we can get consistency, but not \sqrt{n} -consistency or asymptotic normality. However, the fact that the function drops vertically indicates that we may have an “unbounded Hessian,” which could indicate convergence of faster than \sqrt{n} .

The true rate of convergence is actually n . To see this, we'll approximate the CDF of $n(\hat{\alpha} - \alpha_0)$. It is

¹³Actually, in modern terms, *semiparametric*

obviously 1 for non-negative arguments. For arguments $t \leq 0$, we have:

$$\begin{aligned}
\mathbb{P}\{n(\hat{\alpha} - \alpha_0) \leq t\} &= \mathbb{P}\{\max_i X_i \leq \alpha_0 + t/n\} \\
&= \mathbb{P}\{X_1 \leq \alpha_0 + t/n, \dots, X_n \leq \alpha_0 + t/n\} \\
&= \mathbb{P}\{X_1/\alpha_0 \leq 1 + t/(n\alpha_0), \dots, X_n/\alpha_0 \leq 1 + t/(n\alpha_0)\} \\
&= \left(1 + \frac{t}{n\alpha_0}\right)^n \\
&\rightarrow e^{t/\alpha_0}
\end{aligned}$$

Note that this is 1 at $t = 0$ as expected. Thus, $n(\hat{\alpha} - \alpha_0) = O_P(1)$, so the estimator is consistent at rate n , and converges to an exponential distribution.

8 Bootstrapping

Remark. We will not introduce any new estimators, but will talk a lot about asymptotics here. Bootstrapping comes from [Efron \(1979\)](#). It is somewhat magical to statisticians. At a big picture level, we are essentially using the analog principle for inference rather than estimation.

Example. Consider any quantity of interest that can be defined as $\theta = g(F)$ where $g(\cdot)$ is a known function and F is the true distribution of the data. Obvious examples are the mean $\mu = \mathbb{E}X$ or the linear projection $\hat{\beta} = (\mathbb{E}X'X)^{-1}\mathbb{E}X'Y$. Imagine we have an estimator \hat{F} of F . It would be natural to estimate $g(F)$ with $g(\hat{F})$. If \hat{F} is consistent in a sufficiently strong sense and g is continuous, the estimator will be consistent.

In principle, $g(\cdot)$ could also be the standard error or the CDF of a given test statistic at a certain sample size n . Then we could use a plug-in estimator to estimate a standard error or a critical value, which is the basic idea of bootstrap inference. We will assume here that the data are i.i.d., but bootstrapping has been extended beyond that case.

Remark. We will think about two estimands. One is the (scaled) standard deviation of an estimator. The other is a general distribution of a sample statistic

$$J_n(t, F) = \mathbb{P}\{T_n \leq t \mid F\} \quad \text{where} \quad T_n(W_1, \dots, W_n, F) \text{ is a sample test statistic}$$

Hence, F is the population distribution of observables. Note that both T_n and J_n are indexed by sample size n . The idea is to estimate J_n with a plug-in estimator using \hat{F} . In the conventional notation, we write:

$$J_n^* = J_n(t, \hat{F})$$

The simple nonparametric bootstrap estimates F with the empirical distribution F_n , which is not essential and often not optimal, but easy to conceptualize. Formally:

$$\begin{aligned}
J_n^*(t) &= J_n(t, F_n) \\
F_n(w) &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{W_i \leq w\}
\end{aligned}$$

We will estimate the distribution of $\hat{\theta}$ by computing its distribution under (i.i.d. for now) sampling from F_n at sample size n *with replacement*. By the Glivenko-Cantelli Theorem, $F_n(w) - F(w) \xrightarrow{a.s.} 0$ uniformly over $w \in \mathbb{R}^k$. We hope that this makes J_n^* a good estimator of J_n , but that's not so clear. We will later specify conditions that are necessary for it to hold.

Example. Let $\theta = \bar{X}$, $\theta = \mathbb{E}X$, the sample is $\{0, 1, 2\}$, for $n = 3$. To compute the bootstrapped distribution of the sample average, we essentially put the three observations into an urn, sample with replacement, and

increase n . We compute the bootstrapped sample average \bar{X}^* by taking i.i.d. draws from $\{0, 1, 2\}$. The bootstrap distribution from a size 3 resample is:

\bar{X}^*	0	1/2	2/3	1	4/3	5/3	2
p.m.f.	1/27	3/27	6/27	7/27	6/27	3/27	1/27
$(X^* - \bar{X})$	-1	-2/3	-1/3	0	1/3	2/3	1

It should be clear that this is approaching a normal distribution!

Remark. This is, of course, on the order of n^m realizations. In practice, this is not even a little bit possible to compute, basically always it's done through simulation studies – Monte Carlo bootstraps.

Question. What can we do with this? The two main applications are (i) bootstrapped standard errors and (ii) tests and confidence intervals that directly use the bootstrapped CDF of the test statistic, called *percentile methods*.

Remark. Even though bootstraps *feel* really finite sample, all of the justifications we'll cover work only if central limit theorems hold, and the method is basically entirely justified by asymptotics.

There is a very common misleading intuition is that this works because $F_n \rightarrow F$ implies that $F_n(\bar{X}^* - \bar{X}) \approx F(\bar{X} - \mathbb{E}X)$ as n increases. However, though that is true, it's trivial because both converge to a degenerate point mass at 0. That establishes consistency, but not that any confidence intervals are useful in any way. If we estimate the stable asymptotic distribution, so generally $T_n = \sqrt{n}(\bar{X} - \mathbb{E}X)$ and $J_n(t) = \mathbb{P}\{T_n \leq t\}$, which we estimate by $T_n^* = \sqrt{n}(\bar{X}^* - \bar{X})$ and $J_n^*(t) = \mathbb{P}\{T_n^* \leq t\}$. It is completely not obvious that this works! This converges at rate \sqrt{n} , but there is also error at rate \sqrt{n} , so this is $O_P(1)$, which is bad! In fact, this works if and only if the distribution is asymptotically normal.

Question. How would you prove any of this?

	n	∞
Population	$J_n(t)$	$J_\infty(t)$
Bootstrap	$J_n^*(t)$	$J_\infty^*(t)$

We can almost always prove convergence along the rows, so $J_n(t) \rightarrow J_\infty(t)$ and $J_n^*(t) \rightarrow J_\infty^*(t)$. In relatively benign cases, we can also show that $J_\infty^*(t) \rightarrow J_\infty(t)$. A big necessary condition is that the limiting distribution of the test statistic must be continuous in underlying parameters, which is not an innocuous assumption.

Remark. We can divide estimators and test statistics into the following classes:

1. *Pivot*, meaning that $J_n(t, F)$ does not depend on F
2. *Asymptotic Pivot*, meaning that $J_\infty(t, F)$ does not depend on F
3. *Asymptotically Continuous*, meaning that $J_\infty(t, F)$ is continuous in F
4. None of the above

Example. The most common use for bootstrapping is to get confidence intervals. We will explore the two main methods used for this:

1. *Bootstrapped Standard Errors*. We can get precisely

$$SE^* = \left(\mathbb{E}^* \left(\hat{\theta}^* - \mathbb{E}^* \hat{\theta}^* \right)^2 \right)^{1/2}$$

where \mathbb{E}^* is with respect to the (bootstrap) distribution of $\hat{\theta}^*$. Note that in cases where we know that $\mathbb{E}^* \hat{\theta}^* = \hat{\theta}$, this simplifies to $\left(\mathbb{E}^* \left(\hat{\theta}^* - \hat{\theta} \right)^2 \right)^{1/2}$. In particular, this is the case when (i) the estimator is unbiased, and (ii) the empirical estimate is the bootstrap population true value; both are true in particular for $\theta = \mathbb{E}X$. We can approximate this to an arbitrary degree of precision by choosing a high

B in

$$SE^{MC} = \left(\frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}^b - \frac{1}{B} \sum_{b=1}^B \hat{\theta}^b \right)^2 \right)^{1/2}$$

which may simplify to

$$SE^{MC} = \left(\frac{1}{B} \sum_{b=1}^B \left(\hat{\theta}^b - \hat{\theta} \right)^2 \right)^{1/2}$$

Many estimators are \sqrt{n} -consistent and asymptotically normal, but with difficult to estimate asymptotic variances. In these cases, it's convenient to bootstrap the standard errors and then report Wald confidence intervals. Note, however, that we require that the statistic we are estimating is not only standard normal but has the relevant moments. This can fail in relevant cases – see IV or TSLS.

2. **Percentile Interval.** This is basically the exact procedure we would naïvely think works for bootstrapping. It almost entirely works the way we would assume. To motivate, think about the ‘oracle’ confidence interval, which is the interval we would use if all population quantiles were known. Letting $T_n = \hat{\theta} - \theta_0$ with exact quantile function q_n , this confidence interval would be

$$CI^{\text{oracle}} = \left[\hat{\theta} - q_n(1 - \alpha/2), \hat{\theta} - q_n(\alpha/2) \right]$$

because

$$\mathbb{P}\{\theta_0 \in CI^{\text{oracle}}\} = \mathbb{P}\left\{ \hat{\theta} - q_n(1 - \alpha/2) \leq \theta_0 \leq \hat{\theta} - q_n(\alpha/2) \right\} = \underbrace{\mathbb{P}\left\{ \hat{\theta} - q_n(\alpha/2) \leq \hat{\theta} - \theta_0 \leq q_n(1 - \alpha/2) \right\}}_{=1-\alpha}$$

The bootstrap percentile interval is just a plug-in estimator for this object, and is extremely easy to compute:

Algorithm. Take α as exogenous.

- (a) Generate B bootstrap realizations $\hat{\theta}^1, \dots, \hat{\theta}^B$. Let the vector $(\hat{\theta}^{[1]}, \dots, \hat{\theta}^{[B]})$ collect them in increasing order.
- (b) Up to integer constraints, the bootstrap percentile $(1 - \alpha)$ confidence interval is

$$\left[2\hat{\theta} - \hat{\theta}^{[(1-\alpha/2)B]}, 2\hat{\theta} - \hat{\theta}^{[\alpha B/2]} \right]$$

9 Non-Parametrics

Remark. In modern usage, the term *non-parametric* does not mean the lack of parameters or even parameterizations. Rather, it refers to the data generating process is specified up to an infinite dimensional unknown quantity. A parametric model is specified up to a finite dimensional parameter. This is confusing, but it's the usage we have.

Remark. We will develop only Kernel Density Estimation here. This is not particularly relevant for empiricists, but almost all actual non- or semi-parametric models build on the theory of kernel density estimation, and it's actually tractable for us to think about.

Model. Say that we want to estimate the distribution of a random variable X_i . An obvious estimator for the CDF F , the empirical distribution, was already discussed and used in the bootstrap. The analogous estimator of the density is the empirical probability mass function, which consists of at most n mass points. This is often useless – the empirical p.m.f. consistently estimates the density almost nowhere.

Our fix to this is *kernel density estimation*, where we take a weighted average of nearby mass points of the empirical p.m.f. This can be intuited as smoothing out histograms; or as estimating F' by some smoothed arc slope of F_n . The basic idea is obvious, but there are many ways to do the smoothing. Specifically, we have to choose a *kernel* and a *bandwidth*. In practice, the former doesn't matter so much but the latter is significant and affects the results a lot.

In the scalar case, a kernel density estimator of $f(x)$ can be written as

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{X_i - x}{h}\right)$$

where $k : \mathbb{R} \rightarrow \mathbb{R}$ is the kernel function and h is the bandwidth.

Example. For a simple example, let $k(t) = \frac{1}{2} \cdot \mathbb{1}\{|t| \leq 1\}$, the *uniform kernel*, and we can write

$$\hat{f}(x) = \frac{1}{2nh} \sum_{i=1}^n \mathbb{1}\{x - h \leq X_i \leq x + h\} = \frac{1}{2h} P_n([x - h, x + h])$$

where P_n is the empirical p.m.f.

Remark. This particular estimator is a step function. This could be avoided by using continuous kernels, but the uniform kernel makes some of the trade-offs in kernel estimation extremely obvious. In particular, we estimate the density at x by averaging the empirical p.m.f. over $[x - h, x + h]$. A larger choice of h means that we average over more observations, so the variance of our estimator will decrease. However, it also means that we average over a set where the sample density may be increasingly different from the true density at x , so the bias will increase. This is similar to histograms: if we choose bins that are too large, the histogram eventually becomes constant. If we choose bins that are too small, it will become too wiggly and will include unwanted zeros.

Remark. In practice, we will choose h to resolve the above issues. In an oracle situation where we know the density to be estimated, we can find the choice of h that minimizes the MSE $\mathbb{E}(\hat{f}(x) - f(x))^2$ at some x of interest. It is intuitive that the optimal h here will (i) decrease with n , and (ii) depend on how smooth the true density is at x .

The obvious issue is that we definitionally do not know the true $f(x)$. This problem is considerable and there is a large literature. Additionally, the optimal choice of h as a function of n will imply that bias and variance are of the same order as n . This is because we want to minimize the larger of two rates, and the maximum of two functions is minimized at a point of equality. So the bandwidth choice that minimizes MSE will also lead to an asymptotic distribution that is non-centered. It is common in practice to choose a smaller bandwidth than optimal so that the variance dominates the bias and the asymptotic distribution is centered, a process called *(deliberate) undersmoothing*.

Definition. A *kernel function* $k : \mathbb{R} \rightarrow \mathbb{R}$ must integrate to 1: $\int k(u) \partial u = 1$. No other property is strictly required, but in practice kernel functions are typically symmetric, so $k(u) = k(-u)$.

The *order* of k is the index of its first nonzero moment. Formally, defining $\kappa_j(k) = \int u^j k(u) \partial u$, the order of k equals $\min\{j \in \mathbb{Z}_+ : \kappa_j(k) \neq 0\}$. Since we restrict attention to symmetric kernels and all odd moments of symmetric kernels are zero, any kernel that we look at will have an even order of at least 2. The most popular kernels are non-negative, so will have order of exactly 2. A kernel is *higher order* if its order strictly exceeds 2. We will occasionally mention properties of higher order kernels, but our focus will be on non-negative kernels.

Remark. Any non-negative kernel can be interpreted as a probability density, and so any kernel density estimate using a non-negative kernel is a weighted average of the empirical p.m.f. Two simple and common

kernels are the *uniform* and *Gaussian* kernels:

$$k^{\text{uni}}(u) = \frac{1}{2} \mathbb{1}\{|u| \leq 1\}$$

$$k^{\text{gau}}(u) = (2\pi)^{-1/2} \exp(-u^2/2)$$

With these (or any other non-negative kernels!) the kernel density estimator can be intuited as taking each of the n realizations of X_i and replacing it with $1/n$ multiplied by a p.d.f. (normal, uniform, etc) centered at that realization. This makes it obvious that $\hat{f}(x)$ is itself a proper p.d.f., but also that the random variable described by \hat{f} is a mean-preserving spread of the random variable described by the empirical distribution.¹⁴

It's intuitively obvious, but we can also see formally why $\hat{f}(x)$ integrates to 1. Note first that

$$1 = \int k(u) \partial u = \int \frac{1}{h} k\left(\frac{X_i - x}{h}\right) \partial x$$

where the last step uses the change of variable $u(x) = \frac{X_i - x}{h}$. Next, we have that

$$\int \hat{f}(x) \partial x = \int \frac{1}{nh} \sum_{i=1}^n k\left(\frac{X_i - x}{h}\right) \partial x = \frac{1}{n} \sum_{i=1}^n \int \frac{1}{h} k\left(\frac{X_i - x}{h}\right) \partial x = 1$$

We can attain moments of the kernel density estimator the obvious way.

Remark. The next step is to develop asymptotics (bias, variance, and MSE). The goal is to (i) figure out a good choice of k and h , and (ii) conduct inference. Getting there requires some algebra.

First, start by writing

$$\mathbb{E}\hat{f}(x) = \int \frac{1}{h} k\left(\frac{t - x}{h}\right) f(t) \partial t = \int k(u) f(x + hu) \partial u$$

The intuition here is that \hat{f} essentially averages the estimated density over a neighborhood of order h of x . Note however that this argument was not asymptotic – without approximations, we would be stuck here. We can move forward by taking the Taylor expansion, of order ν :

$$f(x + hu) = f(x) + f'(x)hu + \frac{1}{2}f''(x)h^2u^2 + \cdots + \frac{1}{\nu!}f^{(\nu)}(x)h^\nu u^\nu + o(h^\nu)$$

Remark. We need here to assume that the partials up to $\nu + 1$ exist. This is not just a regularity condition – as we consider higher order kernels, assuming that the density is sufficiently smooth gets very restrictive quickly.

We next use linearity of integrals to get

$$\begin{aligned} \mathbb{E}\hat{f}(x) &= \int k(u)f(x)\partial u + \int k(u)f'(x)h u \partial u + \frac{1}{2} \int k(u)f''(x)h^2 u^2 + \cdots + \frac{1}{\nu!} \int k(u)f^{(\nu)}(x)h^\nu u^\nu \partial u + o(h^\nu) \\ &= f(x) + f'(x)h \int k(u)u \partial u + \frac{f''(x)h^2}{2} \int k(u)u^2 + \cdots + \frac{f^{(\nu)}(x)h^\nu}{\nu!} \int k(u)u^\nu \partial u + o(h^\nu) \\ &= f(x) + \frac{1}{2}f''(x)h^2 \kappa_2 + \cdots + \frac{1}{\nu!}f^{(\nu)}(x)h^\nu \kappa_\nu + o(h^\nu) \\ &= f(x) + \frac{1}{\nu!}f^{(\nu)}(x)h^\nu \kappa_\nu + o(h^\nu) \end{aligned}$$

because the kernel is of order ν . This makes it clear why we might want to use higher order kernels in theory.

¹⁴In contrast, with higher order kernels, $\hat{f}(x)$ can take negative values and therefore need not be a density; however it generally replicates the second (or higher) moment of the empirical distribution.

In practice, with non-negative kernels we always have

$$\mathbb{E}\hat{f}(x) = f(x) + \frac{1}{2}f''(x)h^2\kappa_2 + O(h^4) \implies \text{bias}(\hat{f}(x)) = \frac{1}{2}f''(x)h^2\kappa_2 + O(h^4)$$

so the bias is always of order $O(h^2)$.

Similarly, we can take the variance:

$$\begin{aligned} \text{Var}(\hat{f}(x)) &= \text{Var}\left(\frac{1}{nh}\sum_{i=1}^n k\left(\frac{X_i - x}{h}\right)\right) \\ &= \frac{1}{nh^2}\text{Var}\left(k\left(\frac{X_i - x}{h}\right)\right) \\ &= \frac{1}{nh^2}\mathbb{E}\left(k\left(\frac{X_i - x}{h}\right)^2\right) - O(1/n) \\ (\text{Taylor}) &\approx \frac{1}{nh}f(x)R(k) + O(1/n) \end{aligned}$$

where the last step (and a bunch of algebra we skipped) defined the *roughness* $R(g) = \int g(t)^2 \partial t$ of a function $g : \mathbb{R} \rightarrow \mathbb{R}$.

We can find the MSE of a certain estimator the standard way:

$$\text{MSE}[\hat{f}(x)] = \left(\text{bias}[\hat{f}(x)]\right)^2 + \text{Var}[\hat{f}(x)] = \left(\frac{1}{\nu!}f^{(\nu)}(x)h^\nu\kappa_\nu + o(h^\nu)\right)^2 + \frac{1}{nh}f(x)R(k) + O(1/n)$$

We will always choose h such that in the above, all terms are dominated other than the *asymptotic mean square error (AMSE)*:

$$\text{AMSE}[\hat{f}(x)] = \left(\frac{1}{\nu!}f^{(\nu)}(x)h^\nu\kappa_\nu\right)^2 + \frac{1}{nh}f(x)R(k)$$

However, which of the two terms in this expression dominates depends on our choice of h – this is precisely the bias–variance tradeoff!

We will choose h and n such that $h \rightarrow 0$ but $nh \rightarrow \infty$. These requirements make sense intuitively – we need bandwidth to vanish for bias to vanish, but it must vanish slowly enough that in expectation the sample size effectively used diverges rather than degenerating.

Remark. We could try to specify k and h to optimize $\text{AMSE}[\hat{f}(x)]$. However, the solution would typically depend on x . If we are interested in the overall performance of the estimator, we usually integrate to get the *asymptotic mean integrated square error (AMISE)*:

$$\text{AMISE}[\hat{f}] = \int_{-\infty}^{\infty} \left[\left(\frac{1}{\nu!}f^{(\nu)}(x)h^\nu\kappa_\nu\right)^2 + \frac{1}{nh}f(x)R(k) \right] \partial x = \left(\frac{1}{\nu!}h^\nu\kappa_\nu\right)^2 R(f^{(\nu)}) + \frac{R(k)}{nh}$$

So in the salient case of non-negative kernels, we have that

$$\text{AMISE}[\hat{f}] = \frac{1}{4}h^4\kappa_2^2R(f'') + \frac{R(k)}{nh}$$

Remark. If we use the optimal bandwidth, bias and variance will be of the same order. This implies something like:

$$\sqrt{nh_n^*}(\hat{f}(x) - f(x)) \xrightarrow{d} \mathcal{N}(\text{abias}, \text{aVar})$$

However, as we saw above, the bias and the standard deviation are of the same order, $\sqrt{nh_n^*}$. This means that the asymptotic bias will *not* go to zero – this will be a non-centered normal. Note also that we do not explicitly know the asymptotic bias unless we have a closed form for f , which in practice is naturally impossible. How do we get inference here?

In practice, we do so by choosing a bandwidth that is *too small* – we take advantage of the bias–variance tradeoff by deliberately *undersmoothing*! We allow ourselves to have a higher variance to get a smaller bias. This will allow the bias to degenerate to zero, and we can get inference!

Remark. There are two downsides to this:

1. Since the confidence interval can always be interpreted as a set-valued estimator, this now converges slower than the optimal point estimator. In large samples, this is extremely large compared to the optimal estimator \pm two standard errors, and is not guaranteed to cover the optimal estimator. We could report this along with the corresponding estimator, but we would sacrifice precision (and communication)
2. While we claim we are degenerating the bandwidth at a certain rate, in practice we of course choose a fixed bandwidth. In principle, we could claim any bandwidth is part of a sequence that converges at our rate of choice. There’s no particularly principled way to do this, which (contrary to what you would think!) empiricists dislike, as one could always claim they chose a suboptimal bandwidth to overestimate results, or other such things.

Question. How should we choose bandwidth in practice? The optimal bandwidth is a complicated function of unknown functions. In practice, we either perform *cross-validation* or we *pre-estimate* in a parametric class of estimators (e.g. *Silverman’s Rule of Thumb*).

Question. Can we find the optimal kernel? In principle, yes! We intuitively try and find the minimum roughness such that k is still a density. The solution is the *Epanechnikov Kernel*. Also popular is the Gaussian kernel.

Remark. Frankly, any reasonable kernel choice will generally be fine in practice. The reason Gaussian kernels are popular is because they’re everywhere infinitely continuously differentiable, which the Epanechnikov kernel is not. The Epanechnikov kernel is the default in Stata and R, however.

In practice, bandwidth choice is a much bigger issue than kernel choice. Changing the kernel won’t change the results much, but bandwidth choice changes a ton.

Remark. Finally, kernel density doesn’t work in many dimensions. As you increase the dimensionality, the rate of convergence gets slower and slower, and indeed actually converges slowly enough that it fully does not work. This is the *curse of dimensionality*. The optimal rate of convergence in general will be $O(n^{-4/(4+D)})$. The highest dimension that works in practice is 2.