

# EMPIRICAL ANALYSIS II-2.

WINTER 2019

(LARS HANSEN)

SIMON SANGMIN OH

UNIVERSITY OF CHICAGO

**Note:**

These supplementary notes are meant to be used in conjunction with the lecture notes provided by Prof. Hansen.

## Contents

<b>1</b>	<b>Decision Theory</b>	<b>3</b>
1.1	Bayesian Decision Rules . . . . .	3
1.2	Partial Orderings and Admissibility . . . . .	6
1.3	Max-min Decision Rules . . . . .	6
1.4	Risk vs. Uncertainty . . . . .	13
<b>2</b>	<b>Stochastic Processes</b>	<b>14</b>
2.1	Setup . . . . .	14
2.2	Law of Large Numbers . . . . .	14
2.3	Limiting Empirical Measures . . . . .	15
2.4	Ergodic Decomposition . . . . .	15
2.5	Example . . . . .	16
<b>3</b>	<b>Markov Processes</b>	<b>17</b>
3.1	Finite-state Markov Chain . . . . .	17
<b>4</b>	<b>Additive Functionals</b>	<b>18</b>
4.1	Setup . . . . .	18
4.2	Central Limit Theory . . . . .	18
4.3	Martingale Decomposition Galore . . . . .	20
4.4	TA Session: Permanent vs Transient Shocks . . . . .	23
4.5	TA Session: Small Shock Approximation a la Lombardo and Uhlig (2018) . . . . .	24
<b>5</b>	<b>Likelihood Processes (Ch. 8)</b>	<b>28</b>
5.1	Likelihood Constructions . . . . .	28
5.2	Likelihood Ratios . . . . .	28
5.3	Score Processes . . . . .	30
5.4	Nuisance Parameters . . . . .	31
<b>6</b>	<b>Learning (Ch. 9)</b>	<b>33</b>
6.1	Learning about discrete states . . . . .	33
<b>7</b>	<b>Generalized Method of Moments</b>	<b>36</b>
7.1	Preliminaries . . . . .	36
7.2	Traditional Way of Presenting GMM (via Minimization Problem) . . . . .	38
7.3	GMM Estimation with Constant Selection . . . . .	38
7.4	GMM Limiting Approximation . . . . .	40
7.5	GMM Efficiency Bound . . . . .	41
7.6	Testing with GMM . . . . .	44
7.7	Arguing for Consistency . . . . .	46
7.8	Application: Exchange Rates . . . . .	46
7.9	Key GMM Equations . . . . .	46
7.10	Revisiting Linear Models via GMM . . . . .	47

# 1 Decision Theory

## 1.1 Bayesian Decision Rules

We consider the problem of a decision maker having to make a decision after observing a realization of a random vector, the distribution of which is not completely known.

### 1.1.1 Basic Setup

Consider the following setup:

- ▷ A random vector  $Y \in \mathcal{Y}$  described by a probability density  $\psi(y|\theta)$  relative to a measure  $\tau$  over  $\mathcal{Y}$ . The likelihood function  $\psi(y|\theta)$  can be thought of as a *statistical model*.
  - ▷  $\theta \in \Theta$  is an unknown parameter that affects the probability distribution of  $Y$ .
  - ▷ Decision rule  $D : \mathcal{Y} \rightarrow \mathcal{D}$  (Borel measurable).
  - ▷ Preferences represented by a utility function  $U(d, \theta, y)$ , where  $y \in \mathcal{Y}$  denotes a realisation of  $Y$ , and  $d \in \mathcal{D}$  is the decision made by the DM having observed  $y$  (it can depend on  $y$ ).
- \* The  $\theta$  term may show up in the utility function because the  $\theta$ s, in a dynamic context, inform the forecast of the consequences of my action.

I want to adjust this utility for risk. In other words, to consider the  $U$  before observing  $y$ , we can compute the expectation and define the *risk function* as

$$\bar{U}(D|\theta) = \int_{\mathcal{Y}} U[D(y), \theta, y] \psi(y|\theta) \tau(dy).$$

$\bar{U}$  is referred to as a *risk function* (in statistical decision theory) and it depends on both the decision rule  $D$  and the parameter  $\theta$ . Integration of  $U$  over  $y$  conditional on  $\theta$  adjusts for the risk in  $y$ . This is analogous to the “frequentist risk” discussed in Prof. Uhlig’s class.

**Example 1.1.** (*Model selection problem*). Suppose that  $\Theta = \{\theta_1, \theta_2\}$ , where each value of  $\theta$  corresponds to a statical model (i.e. a different probability density over  $\mathcal{Y}$ ). A decision-maker wishes to select a model and has the following utility function:

$$U(d, \theta, y) = \begin{cases} d & \text{if } \theta = \theta_1 \\ 1 - d & \text{if } \theta = \theta_2 \end{cases},$$

where  $d \in \mathcal{D} = [0, 1]$  is the probability of choosing  $\theta_1$  and  $1 - d$  is the probability of choosing  $\theta_2$  (this specification allows for randomised decisions).

We focus on threshold decision rules. Partition  $\mathcal{Y}$  into two disjoint sets  $\mathcal{Y}_1$  and  $\mathcal{Y}_2$ , and consider the non-randomised decision rule:

$$D(y) = \begin{cases} 1 & \text{if } y \in \mathcal{Y}_1 \\ 0 & \text{if } y \in \mathcal{Y}_2 \end{cases}.$$

The risk function then for an arbitrary partitioning is given as:

$$\begin{aligned}
 \bar{U}(D|\theta_1) &= \int_{\mathcal{Y}} U[D(y), \theta_1, y] \psi(y|\theta_1) \tau(dy) \\
 &= \int_{\mathcal{Y}_1} U[1, \theta_1, y] \psi(y|\theta_1) \tau(dy) + \int_{\mathcal{Y}_2} U[0, \theta_1, y] \psi(y|\theta_1) \tau(dy) \\
 &= \int_{\mathcal{Y}_1} \psi(y|\theta_1) \tau(dy) = \mathbb{P}(y \in \mathcal{Y}_1|\theta_1), \\
 \bar{U}(D|\theta_2) &= \int_{\mathcal{Y}_1} U[1 - 1, \theta_2, y] \psi(y|\theta_2) \tau(dy) + \int_{\mathcal{Y}_2} U[1 - 0, \theta_2, y] \psi(y|\theta_2) \tau(dy) \\
 &= \int_{\mathcal{Y}_2} \psi(y|\theta_2) \tau(dy) = \mathbb{P}(y \in \mathcal{Y}_2|\theta_2).
 \end{aligned}$$

$\bar{U}(D|\theta_i)$  gives the probability that the decision rule  $D(y)$  “correctly” chooses the model  $\theta_i$ , while  $1 - \bar{U}(D|\theta_i)$  is the probability of choosing the incorrect model.

### 1.1.2 Ex-ante vs. Ex-post Problems

Let a probability measure  $\pi$  over  $\Omega$  denote the *prior distribution* that summarises the DM’s ignorance about  $\theta$ .

**Definition 1.1.** (*Ex ante decision problem*) A Bayesian DM maximises his expected utility *using his prior beliefs* with respect to his decision rule  $D$ ; i.e.

$$\max_D \int_{\Theta} \bar{U}(D|\theta) \pi(d\theta)$$

The solution to the problem is the optimal decision rule  $D^*(y)$ .

**Definition 1.2.** (*Ex post decision problem*) For all  $y \in \mathcal{Y}$ ,

$$\max_{d \in \mathcal{D}} \int_{\Theta} U(d, \theta, y) \bar{\pi}(d\theta|y) := \mathbb{E}[U(d, \theta, y) | y]$$

The maximiser  $d^*$  for each  $y \in \mathcal{Y}$  induces the decision rule  $\bar{D}^*(y)$ .

The two problems differ in the following respects:

1. Observation of  $y \in \mathcal{Y}$ 
  - ▷ The ex-ante problem gets its name because the agent needs to derive the decision rule before he observes  $y \in \mathcal{Y}$ .
  - ▷ The ex-post problem gets its name because the agent now chooses the decision rule after observing  $y \in \mathcal{Y}$ .
2.  $U$  vs.  $\bar{U}$ 
  - ▷ In the ex-ante problem, the only probability measure he can use is his prior, and he will choose a decision rule that maximizes his expected *adjusted* utility ( $\bar{U}$ ) with respect to the prior.
    - \* Since  $y$  is uncertain, we need to adjust  $U$  for the different values of  $y$  and obtain  $\bar{U}$  first.
  - ▷ In the ex-post problem, the agent will use the updated posterior belief that maximizes his expected utility ( $U$ ) with respect to the posterior.
    - \* Here we already observe  $y$ , so we can just use  $U$  in our maximization problem.

### 3. Decomposing ex-ante problems

- ▷ One can break the ex-ante decision problem into two subproblems: (1) Given an observation of  $Y$ , estimate  $\theta$  and (2) Given an estimate of  $\theta$ , choose a decision rule.

Notice that if the objective function of the ex ante decision problem is finite when evaluated at  $D^*$ , we have then  $D^* \equiv \bar{D}^*$ . This result may not hold, however, in an experiment setting where your decisions affect outcomes and the data you generate.

**Example 1.2.** (*Model Selection Revisited*) What does the Bayesian decision rule look like? We start by computing the posterior probabilities:

$$\begin{aligned}\bar{\pi}(d\theta_1|y) &= \frac{\psi(y|\theta_1)\pi_1}{\psi(y|\theta_1)\pi_1 + \psi(y|\theta_2)\pi_2}, \\ \bar{\pi}(d\theta_2|y) &= \frac{\psi(y|\theta_2)\pi_2}{\psi(y|\theta_1)\pi_1 + \psi(y|\theta_2)\pi_2}.\end{aligned}$$

Using these probabilities, for any  $y \in \mathcal{Y}$ , the ex-post decision problem is then:

$$\begin{aligned}\int_{\Theta} U(d, \theta, y) \bar{\pi}(d\theta|y) &= \int_{\{\theta_1, \theta_2\}} U(d, \theta, y) \bar{\pi}(d\theta|y) \\ &= U(d, \theta_1, y) \bar{\pi}(d\theta_1|y) + U(d, \theta_2, y) \bar{\pi}(d\theta_2|y) \\ &= d\bar{\pi}(d\theta_1|y) + (1-d)\bar{\pi}(d\theta_2|y).\end{aligned}$$

and thus:

$$\max_{d \in \mathcal{D}} d\bar{\pi}(d\theta_1|y) + (1-d)\bar{\pi}(d\theta_2|y)$$

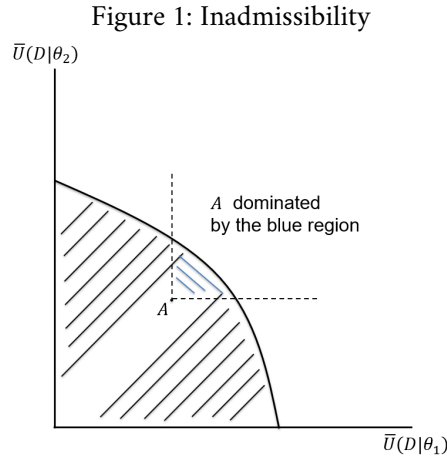
The solution to this problem is then

$$d^*(y) = \begin{cases} 1 & \bar{\pi}(d\theta_1|y) > \bar{\pi}(d\theta_2|y) \\ [0, 1] & \bar{\pi}(d\theta_1|y) = \bar{\pi}(d\theta_2|y), \forall y \in \mathcal{Y}. \\ 0 & \bar{\pi}(d\theta_1|y) < \bar{\pi}(d\theta_2|y) \end{cases}$$

Assuming that  $d^*(y) = 0$  if  $\bar{\pi}(d\theta_1|y) = \bar{\pi}(d\theta_2|y)$ , we can write a *threshold rule* that is a solution to the ex post problem:

$$d^*(y) = \begin{cases} 1 & \bar{\pi}(d\theta_1|y) > \bar{\pi}(d\theta_2|y), \forall y \in \mathcal{Y}. \\ 0 & \bar{\pi}(d\theta_1|y) \leq \bar{\pi}(d\theta_2|y) \end{cases}$$

Thus, the Bayesian solution of the model selection problem is to select the model for which the posterior probability is highest, and this generates a partition of the set  $\mathcal{Y}$ . Furthermore, it can be shown that this will lie on the boundary of the convex set we drew earlier.



## 1.2 Partial Orderings and Admissibility

Earlier, we saw that the Bayesian decision rule  $D^*$  depends on the prior distribution  $\pi$ . However, even before specifying a prior, the utility function  $\bar{U}(D|\theta)$  allows us to identify some decision rules that cannot be dominated.

**Definition 1.3.** Consider decision rules  $D_1 : \mathcal{Y} \rightarrow \mathcal{D}$  and  $D_2 : \mathcal{Y} \rightarrow \mathcal{D}$ . Under a (partial) ordering  $\succsim$ ,  $D_2$  is preferred over  $D_1$ , denoted  $D_2 \succsim D_1$ , if and only if

$$\bar{U}(D_2|\theta) \geq \bar{U}(D_1|\theta), \forall \theta \in \Theta.$$

The ordering is partial because  $\succsim$  ranks some, but not all, pairs of decision rules. For example, it might be that for some  $\theta_1, \theta_2 \in \Theta$ , we have

$$\bar{U}(D_2|\theta_2) > \bar{U}(D_1|\theta_2), \bar{U}(D_1|\theta_1) < \bar{U}(D_2|\theta_1).$$

We restrict our attention to such undominated decision rules, called *admissible* decision rules. We rule out decision rules that gives strictly lower utility under all priors.

**Definition 1.4.** (*Admissibility*) A decision rule  $D$  is *admissible* if there exists no decision rule  $\tilde{D}$  for which  $\tilde{D} \succ D$  with  $\bar{U}(\tilde{D}, \tilde{\theta}) > \bar{U}(D, \tilde{\theta})$  for some  $\tilde{\theta} \in \Theta$ .

*Remark 1.1.* If you want to get an admissible decision rule, you solve the Bayesian problem. You repeat the process for a set of priors, and you get a whole frontier of admissible decision rules as outlined earlier. This still doesn't address the question of where the prior comes from, but at least we have a boundary.

## 1.3 Max-min Decision Rules

These rules are designed to address the uncertainty regarding prior selection. Specifically, a max-min DM wishes to construct a robust rule (with respect to his uncertainty about the prior). He does so by exploring the consequences of a decision rule across alternative priors by formulating a max-min problem that expresses the DM's aversion to his uncertainty about the prior. Note that the chosen prior depends upon the DM's utility function.

Here, we assume that the DM has *multiple* prior probability distributions over  $\theta$ . Specifically, let  $C$  denote a positive convex function of possible prior probability distributions  $\pi$  over  $\theta$ . The function  $C$  reflects the DM's *ambiguity* about the prior  $\pi$ . This is separate from the function  $U$  which expresses DM's risk aversion (aversion to uncertainty over the outcome  $y$  and  $\theta$ ).

**Definition 1.5.** (*Max-min ex ante problem*)

$$\max_D \min_{\pi} \left[ \int_{\Theta} \bar{U}(D|\theta) \pi(d\theta) + C(\pi) \right].$$

Given a decision rule, solve the minimization problem. Rank the decision rules based on the outcomes of the minimization problem and choose the best one. Note that you can formulate it as max-max. We're just implicitly assuming that the agent is risk-averse and is hence interested in problems of the above sort.

The following are examples of a commonly used  $C$ .

1. Denote  $\Pi$  as the set of possible priors and consider

$$C(\pi) = \begin{cases} 0 & \pi \in \Pi \\ \infty & \pi \notin \Pi \end{cases}.$$

This ensures that  $\pi$  being chosen by the max-min DM belongs in the set  $\Pi$ .

2. Suppose that  $\Theta$  has  $n$  elements and that  $\pi^o$  denotes a "reference" prior. Define

$$C(\pi) = \kappa \sum_{j=1}^n \pi_j (\log \pi_j - \log \pi_j^o)$$

which is a convex function in the probabilities.

- ▷ The summation term is referred to as the *relative entropy* of the  $\pi$  distribution with respect to the  $\pi^o$  distribution. Entropy is the expected value of the log-likelihood ratio, where the expectation is taken with respect to  $\pi$ .
- ▷ The function  $C$  is convex and attains minimum at  $\pi = \pi^o$ . The parameter  $\kappa$  is the penalisation factor—a higher value means greater penalty from deviating from the baseline prior distribution  $\pi^o$ .
- ▷ Setting  $\kappa = \infty$  reduces the max-min ex ante problem to the Bayesian (ex ante) problem with a (unique) prior  $\pi^o$ .
- ▷ To verify that  $C(\pi) \geq 0$ , notice that  $r \log r \geq r - 1$  since the convex function always lies above the gradient line. Then:

$$C(\pi) = \kappa \sum_{j=1}^n \frac{\pi_j}{\pi_j^o} \left( \log \frac{\pi_j}{\pi_j^o} \right) \pi_j^o \geq \kappa \sum_{j=1}^n \left( \frac{\pi_j}{\pi_j^o} - 1 \right) \pi_j^o = 0$$

### 1.3.1 Example: Exponential Tilting

Let us solve the max-min ex ante problem assuming the entropy as the  $C$  function (choice #2 from above)

1. We first solve the inner minimization problem, fixing  $D$ , where you are choosing a probability distribution  $\pi$ :

$$\min_{\pi} \sum_{j=1}^n \pi_j [\bar{U}(D|\theta_j) + \kappa (\log \pi_j - \log \pi_j^o)] \quad s.t. \quad \sum_{j=1}^n \pi_j = 1, \pi_j \geq 0, \forall j, .$$

▷ The Lagrangian is given by

$$\mathcal{L} = \sum_{j=1}^n \pi_j [\bar{U}(D|\theta_j) + \kappa (\log \pi_j - \log \pi_j^o)] + \lambda \left( 1 - \sum_{j=1}^n \pi_j \right).$$

▷ The first-order condition with respect to  $\pi_j$  is given by

$$\bar{U}(D|\theta_j) + \kappa (\log \pi_j^* - \log \pi_j^o) + \kappa - \lambda = 0,$$

which can be rearranged to

$$\begin{aligned} \log \pi_j^* &= \log \pi_j^o - \frac{1}{\kappa} \bar{U}(D|\theta_j) + \left( \frac{\lambda}{\kappa} - 1 \right) \\ \Rightarrow \pi_j^* &= \pi_j^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_j) \right) \exp \left( \frac{\lambda}{\kappa} - 1 \right). \end{aligned}$$

▷ Since  $\pi_j^*$  is a probability, they must sum to one:

$$\pi_j^* = \frac{\pi_j^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_j) \right)}{\sum_{i=1}^n \pi_i^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_i) \right)}$$

▷ *Interpretation:*  $\pi^*$  is the optimal prior. Since  $\exp(-x)$  is decreasing in  $x$ ,  $\pi_j^*$  (exponentially) *tilts* probabilities away from reference prior  $\pi_j^o$  towards parameter values  $\theta_j$  with lower conditional expected utilities. A higher  $\kappa$  tends to equalise probabilities across possible values of  $\theta_j$  (if  $\kappa = \infty$ , we have  $\pi_j^* = \pi_j^o / \sum_{i=1}^n \pi_j^o$ ). If  $\bar{U}$  is small, I will twist the original reference probability more towards the parameter that gives rise to that  $\bar{U}$ .

2. Substituting into the objective function, we obtain the minimised objective function:

$$\begin{aligned} & \sum_{j=1}^n \pi_j [\bar{U}(D|\theta) + \kappa (\log \pi_j - \log \pi_j^o)] \\ &= \sum_{j=1}^n \pi_j \left[ \bar{U}(D|\theta) + \kappa \left( \log \frac{\pi_j^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_j) \right)}{\sum_{i=1}^n \pi_i^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_i) \right)} - \log \pi_j^o \right) \right] \\ &= \sum_{j=1}^n \pi_j \left[ \bar{U}(D|\theta) + \kappa \left( -\frac{1}{\kappa} \bar{U}(D|\theta_j) - \log \left( \sum_{i=1}^n \pi_i^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_i) \right) \right) \right) \right] \\ &= -\kappa \log \left( \sum_{i=1}^n \pi_i^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_i) \right) \right) \sum_{j=1}^n \pi_j \\ &= -\kappa \log \left( \sum_{i=1}^n \pi_i^o \exp \left( -\frac{1}{\kappa} \bar{U}(D|\theta_i) \right) \right). \end{aligned}$$

So you're making this adjustment based on the decision problem. Think of this as an “extra curvature adjustments.”

3. Notice that our objective function satisfies:

$$\begin{aligned} & \sum_{j=1}^n \pi_j [\bar{U}(D|\theta) + \kappa (\log \pi_j - \log \pi_j^o)] \\ & \geq \sum_{j=1}^n \pi_j^* \bar{U}(D|\theta) + \kappa \sum_{j=1}^n \pi_j^* (\log \pi_j^* - \log \pi_j^o) \end{aligned}$$



▷ Now suppose

$$\kappa \sum_{j=1}^n \pi_j (\log \pi_j - \log \pi_j^o) \leq \kappa \sum_{j=1}^n \pi_j^* (\log \pi_j^* - \log \pi_j^o)$$

Then it is necessarily the case that

$$\sum_{j=1}^n \pi_j \bar{U}(D|\theta) \geq \sum_{j=1}^n \pi_j^* \bar{U}(D|\theta)$$

This means that as long as you're looking at  $\pi_j'$ s that is less than or equal to the optimal  $\pi_j^*$  for some discrepancy measure (the first inequality), you are guaranteed a higher expected utility (the second inequality). This is a form of a “robustness bound” since you're guaranteed a certain level of expected utility if the first condition is satisfied.

### 1.3.2 Min-max Decision Rules

Whenever you see a max-min problem, you can transform into a two-person zero-sum game. Specifically, we will explore the consequences of changing the order of the minimization and the maximization problem.

**Definition 1.6.** (*Min-max ex post problem*)

$$\min_{\pi \in \Pi} \max_D \int_{\Theta} \bar{U}(D|\theta) \pi(d\theta) + C(\pi).$$

Now we're maximizing for a given  $\pi$  and then later go over all possible values of  $\pi$ . Observe that  $C(\pi)$  appears additively separably in the objective function. Thus, we can first solve

$$\max_D \int_{\Theta} \bar{U}(D|\theta) \pi(d\theta)$$

and reinsert the maximised objective into the min-max problem.

- ▷ This problem is exactly the Bayesian ex ante problem, which we know we can solve via the ex post problem.
- ▷ Denote the Bayesian decision rule as  $D^*$  and assuming that the solution to the minimisation problem,  $\hat{\pi}$ , exists, we can interpret  $\hat{\pi}$  as a “worst-case prior” for the Bayesian decision rule  $D^*$ .
- ▷ Moreover, since  $D^*$  is a Bayesian rule, the max-min decision rule (which coincides with the min-max decision rule under the assumption here) cannot be dominated except possibly on sets to which the worst-case prior assigns zero probability measure.

*Remark 1.2.* This brings us back to prior choice. Savage would have recommended a form of “economic thinking” or “introspection” in choosing a prior, whereas the above result implies a more systematic way of obtaining a “robust prior choice.”

### 1.3.3 Example: Model Selection via Min-Max

Suppose  $\Theta = \{\theta_1, \theta_2\}$  and the density of  $y$  conditional on  $\theta_1, \theta_2$  are  $\psi(y|\theta_1)$  and  $\psi(y|\theta_2)$ . Utilities are

$$U(d, y, \theta_1) = d, \quad U(d, y, \theta_2) = 1 - d \text{ where } d \in [0, 1]$$

Now a Bayesian with prior  $(\pi_1, \pi_2)$  where  $\pi_1$  is the prior for model 1 to be the correct model and  $\pi_2$  is the prior for model 2 to be the correct model with  $\pi_1 + \pi_2 = 1$ .

### 1. The Inner Maximization Problem:

▷ Taking the ex-post approach:

$$\max_d U(d, y, \theta_1) \bar{\pi}_1(y) + U(d, y, \theta_2) \bar{\pi}_2(y)$$

which is equivalent to

$$\max_d \{ \bar{\pi}_1(y) - \bar{\pi}_2(y) \} d + \bar{\pi}_2(y)$$

which yields the decision rule:

$$D^*(y) = \begin{cases} 1 & \text{if } \bar{\pi}_1(y) > \bar{\pi}_2(y) \\ 0 & \text{if } \bar{\pi}_1(y) \leq \bar{\pi}_2(y) \end{cases}$$

▷ Now we need to find the posterior probabilities:

$$\bar{\pi}_1(y) = \zeta(y) \pi_1 \psi(y|\theta_1)$$

$$\bar{\pi}_2(y) = \zeta(y) \pi_2 \psi(y|\theta_2)$$

where  $\zeta(y)$  is a constant that could potentially depend on  $y$ . Then

$$\begin{aligned} \bar{\pi}_1(y) &> \bar{\pi}_2(y) \\ \Leftrightarrow \pi_1 \psi(y|\theta_1) &> \pi_2 \psi(y|\theta_2) \\ \Leftrightarrow \log \pi_1 + \log \psi(y|\theta_1) &> \log \pi_2 + \log \psi(y|\theta_2) \end{aligned}$$

Denoting  $r \equiv \pi_1/\pi_2$  and  $\phi(y) \equiv \psi(y|\theta_2)/\psi(y|\theta_1)$  we can rewrite the decision rule as

$$D^*(y) = \begin{cases} 1 & \text{if } \phi(y) < r \\ 0 & \text{if } \phi(y) > r \end{cases}$$

▷ Note that  $r$  is a constant that is equal to

$$r = \frac{\pi_1}{1 - \pi_1}$$

because this is a Bayesian's problem who is given the prior  $(\pi_1, \pi_2)$ . There is no min-max or max-min problem here – a nice, simple Bayesian's problem.

▷ We further know that for every  $r \in [0, +\infty]$ , there must exist a prior  $(\pi_1, \pi_2)$  such that

$$D^r(y) = \arg \max_{D(y)} \pi_1 \bar{U}[D|\theta_1] + \pi_2 \bar{U}[D|\theta_2]$$

which is the ex-ante problem for the Bayesian.

▷ Define  $z$  to be the likelihood ratio:

$$z = \phi(y) = \frac{\psi(y|\theta_2)}{\psi(y|\theta_1)}.$$

We can write  $\bar{U}(D|\theta_1)$  and  $\bar{U}(D|\theta_2)$  equivalently as

$$\begin{aligned} \bar{U}(D|\theta_1) &= \int_{\phi(y) < r} \psi(y|\theta_1) \tau(dy), \\ \bar{U}(D|\theta_2) &= \int_{\phi(y) \geq r} \phi(y) \psi(y|\theta_1) \tau(dy). \end{aligned}$$

- ▷ Suppose that  $\psi(y|\theta_1) \tau(dy)$  in  $\mathcal{Y}$  implies a density  $v$  over the positive real numbers in  $z$  space, then

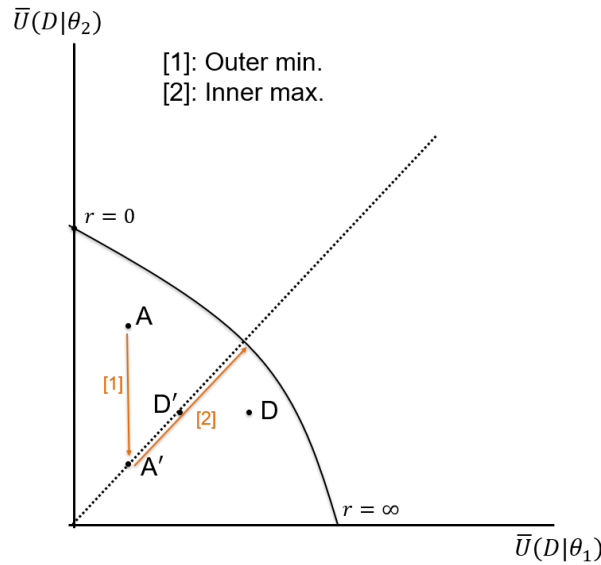
$$\begin{aligned}\bar{U}(D|\theta_1) &= \int_0^r v(z) dz, \\ \bar{U}(D|\theta_2) &= \int_r^\infty zv(z) dz.\end{aligned}$$

Observe that

$$\begin{aligned}\frac{\partial \bar{U}(D|\theta_1)}{\partial r} &= v(r) > 0, & \frac{\partial \bar{U}(D|\theta_2)}{\partial r} &= -rv(r) < 0, \\ \frac{\partial^2 \bar{U}(D|\theta_1)}{\partial r^2} &= v'(r), & \frac{\partial^2 \bar{U}(D|\theta_2)}{\partial r^2} &= -rv'(r) - v(r).\end{aligned}$$

Recall that a minimising DM takes as given the decision rule, which is characterised by the threshold in this case

## 2. Justification of Priors that Satisfy $\bar{U}(D|\theta_1) = \bar{U}(D|\theta_2)$ :



- ▷ Interpretation of the graph:

- \* Each axis is frequentist risk, and each point represents a decision rule. There's nothing Bayesian here – no prior.
  - \* For each decision rule, compute  $\bar{U}(D|\theta_1)$  which is the expected utility assuming model 1 is correct, and  $\bar{U}(D|\theta_2)$  which is the expected utility assuming model 2 is correct.
- ▷ We're going to argue that points off the 45-degree line cannot be the solution to the min-max problem. Suppose that, for a given  $r$  (and therefore  $\pi_1$  and  $\pi_2$ ), it is the case that  $\bar{U}(D|\theta_1) < \bar{U}(D|\theta_2)$ . This is point D in the graph.
- \* Then, the minimising DM (currently solving the outer problem) will put all prior probability on  $\theta_1$  so that minimised objective is given by  $\bar{U}(D|\theta_1)$ .
  - \* Because  $\bar{U}(D|\theta_1)$  is increasing in  $r$ , the maximising DM can increase utility by increasing  $r$ . If he keeps increasing  $r$ , then at some point, we will have  $\bar{U}(D|\theta_1) > \bar{U}(D|\theta_2)$  so that the minimising DM will now wants to put all prior probability on  $\theta_2$  so that the minimised objective is given by  $\bar{U}(D|\theta_2)$ . But since  $\bar{U}(D|\theta_2)$  is decreasing in  $r$ , the maximising DM can increase utility by decreasing  $r$ .

- \* These incentives leads us to search for a value of  $r$  such that

$$\bar{U}(D|\theta_1) = \bar{U}(D|\theta_2).$$

which correspond to points  $A'$  and  $D'$ . At this value of  $r$ , the minimising DM who chooses probabilities has no way to influence the objective.

- \* For the maximising DM, causing  $r$  to deviate from this value would also reduce either side of the equality and therefore diminish the minimised objective function. Thus, the condition above characterises the solution to the max-min problem.
- ▷ As our final solution, the agent will choose the intersection of the frontier and the 45-degree line. This follows directly from admissibility of Bayesian estimators.
- ▷ Another way to see above is to consider the equivalent max-min DM problem.
  - \* First, notice that the solution to the inner minimisation problem for the max-min DM is

$$\min \{ \bar{U}(D|\theta_1), \bar{U}(D|\theta_2) \}.$$

- \* Let  $D^*$  be such that  $\bar{U}(D^*|\theta_1) = \bar{U}(D^*|\theta_2)$ . For any  $D \neq D^*$ , it must be that

$$\bar{U}(D^*|\theta_1) > \bar{U}(D|\theta_1) \text{ or } \bar{U}(D^*|\theta_2) > \bar{U}(D|\theta_2).$$

- \* That is,

$$\min \{ \bar{U}(D^*|\theta_1), \bar{U}(D^*|\theta_2) \} > \min \{ \bar{U}(D|\theta_1), \bar{U}(D|\theta_2) \}$$

so that  $D^*$  solves the max-min DM's maximisation problem, since we want to maximize the min.

### 3. Concavity of the Curve: To simplify notation, let

$$\bar{u}_1(r) = \bar{U}(D|\theta_1), \bar{u}_2(r) = \bar{U}(D|\theta_2).$$

We would like to plot the curve  $(\bar{u}_1(r), \bar{u}_2(r))$  traced out by changing  $r$  over the positive real numbers. In other words, we want to plot

$$\bar{u}_2(r) = f(\bar{u}_1(r)).$$

The slope is given by

$$\begin{aligned} \bar{u}_2'(r) &= f'(\bar{u}_1(r)) \bar{u}_1'(r) \\ \Rightarrow f'(\bar{u}_1(r)) &= \frac{-rv(r)}{v(r)} = -r < 0, \end{aligned}$$

We can also show that it is concave:

$$\begin{aligned} \bar{u}_2''(r) &= f''(\bar{u}_1(r)) (\bar{u}_1'(r))^2 + f'(\bar{u}_1(r)) \bar{u}_1''(r) \\ \Rightarrow -rv'(r) - v(r) &= f''(\bar{u}_1(r)) (v(r))^2 - rv'(r) \\ \Rightarrow f''(\bar{u}_1(r)) &= -\frac{1}{v(r)} < 0. \end{aligned}$$

Hence, the curve is concave.

### 1.3.4 A Modified Example

The curve examined above is symmetric around the 45 degree line given our utility function. What happens if we have the following?:

$$U(d, \theta, y) = \begin{cases} 1.5d & \text{if } \theta = \theta_1 \\ 1 - d & \text{if } \theta = \theta_2 \end{cases}.$$

- ▷ In this case, the intersection with the  $x$ -axis moves to the right, and the intersection with the  $y$ -axis remains the same.
- ▷ To see this, note that the ex post problem is given by

$$\max_{d \in [0,1]} 1.5d\bar{\pi}_1 + (1 - d)\bar{\pi}_2.$$

- ▷ The optimal decision rule (assuming the same tie-break rule):

$$d^*(y) = \begin{cases} 1 & \text{if } 1.5\bar{\pi}(d\theta_1|y) > \bar{\pi}(d\theta_2|y) \\ 0 & \text{if } 1.5\bar{\pi}(d\theta_1|y) \leq \bar{\pi}(d\theta_2|y) \end{cases}, \forall y \in \mathcal{Y}.$$

- ▷ To obtain the max-min solution, we need:

$$\begin{aligned} \bar{U}(d, \theta_1, y) &= 1.5\bar{\pi}(d\theta_1|y) = \bar{\pi}(d\theta_2|y) = U(d, \theta_2, y) \\ &\Rightarrow \frac{\psi(y|\theta_2)}{\psi(y|\theta_1)} = \frac{2}{3} \frac{\pi_1}{1 - \pi_1} = r. \end{aligned}$$

- ▷ We are still going to have the same 45-degree line, but the frontier will look different (larger intercept on the  $x$ -axis).

## 1.4 Risk vs. Uncertainty

Knight (1921) distinguished between risk and uncertainty.  $\psi(y|\theta)$  can be thought of as ex-ante risk; curvature in  $U$  can be thought of as risk aversion. Assigning a prior  $\pi$  over  $\theta$  can be thought of as uncertainty. Below, we discuss one tractable way to capture this notion of uncertainty.

The LHS represents one-stage reduced lottery; the RHS represents a two-stage lottery:

$$\phi^*(y^*|d, y) = \int \underbrace{\phi^*(y^*|d, y, \theta)}_{= \text{risk}} \underbrace{\bar{\pi}^*(\theta|y)}_{= \text{ambiguity}}$$

1. Choose  $y^*$  (next state) given  $\theta$
2. Choose  $\theta$  given current  $\varphi(y^*|y, D, \theta)$

This allows us to separate out the attitude towards risk and the attitude towards ambiguity.

*Remark 1.3.* In finance and macro, you end up having to appeal to having a high risk aversion. The decomposition above allows us to explore a different channel to match the empirical evidence.

*Remark 1.4.* You can never be objective in evaluating empirical evidence. I think that in doing empirical work, you need to make sure where the inputs are coming from. Robustness is one way to show this.

## 2 Stochastic Processes

### 2.1 Setup

- ▷ Invariant event: An event  $\Lambda \in \mathfrak{F}$  is *invariant* if  $\Lambda = \mathbb{S}^{-1}(\Lambda) \cdot \Omega$  and  $\emptyset$  are both invariant events. If  $\Lambda$  is an invariant event, then

$$\mathbb{S}^t(\omega) \in \Lambda, \forall \omega \in \Lambda, \forall t = 0, 1, 2, \dots;$$

i.e. all sample points stay in the event  $\Lambda$  over time.

- ▷ Measure-preserving: The transformation  $\mathbb{S}$  is measure-preserving if

$$\mathbb{P}(\Lambda) = \mathbb{P}\{\mathbb{S}^{-1}(\Lambda)\}$$

for all  $\Lambda \in \mathfrak{F}$  where  $\Lambda$  is an event and  $\mathfrak{F}$  is a collection of events.

- ▷ Conditional expectation: Let  $\{\Lambda_j\}_{j=1}^{\infty}$  denote a countable partition of  $\Omega$ ,<sup>1</sup> where each  $\Lambda_j$  is an invariant event, and  $\mathfrak{J}$  denote the collection of invariant events (a  $\sigma$ -algebra). The expectation conditional on a invariant set is

$$\mathbb{E}[X|\Lambda_i] = \frac{\int_{\Lambda_j} X(\omega) d\mathbb{P}(\omega)}{\mathbb{P}(\Lambda_j)}$$

so you have a different conditional expectation for the corresponding invariant set. The expectation conditional on the collection of all invariant sets  $\mathfrak{J}$  is

$$\mathbb{E}[X|\mathfrak{J}](\omega) = \frac{\int_{\Lambda_j} X(\omega) d\mathbb{P}(\omega)}{\mathbb{P}(\Lambda_j)} \text{ if } \omega \in \Lambda_j.$$

The conditional expectation is constant for  $\omega \in \Lambda_j$  but can vary across  $\Lambda_j$ 's.

- ▷ Ergodicity: A measure-preserving transformation  $\mathbb{S}$  is said to be *ergodic* under a probability measure  $\mathbb{P}$  if all invariant events have probability zero or one.

### 2.2 Law of Large Numbers

- ▷ Law of Large Numbers (*Birkhoff*): Suppose  $\mathbb{S}$  is **measure-preserving**.

1. (Almost-sure convergence) For any  $X$  such that  $\mathbb{E}[|X|] < \infty$ ,

$$\frac{1}{T} \sum_{t=1}^T X_t(\omega) \xrightarrow{\mathbb{P}} \mathbb{E}[X|\mathfrak{J}](\omega).$$

2. (Mean-square convergence) For any  $X$  such that  $\mathbb{E}[|X|^2] < \infty$ ,

$$\mathbb{E} \left[ \left| \frac{1}{T} \sum_{t=1}^T X_t(\omega) - \mathbb{E}[X|\mathfrak{J}](\omega) \right|^2 \right] \rightarrow 0.$$

- ▷ Law of Large Numbers (Standard): (Standard Law of Large Numbers) Suppose  $\mathbb{S}$  is **ergodic** (i.e. for all  $\Lambda \in \mathfrak{F}$ ,  $\mathbb{P}(\Lambda) = 0, 1$ )

<sup>1</sup>That is,  $\Lambda_j \cap \Lambda_k = \emptyset, \forall j \neq k, \bigcup_{j=1}^{\infty} \Lambda_j = \Omega$  and each  $\Lambda_j$  is nonempty.

1. (*Almost-sure convergence*) For any  $X$  such that  $\mathbb{E}[|X|] < \infty$ ,

$$\frac{1}{T} \sum_{t=1}^T X_t(\omega) \xrightarrow{\mathbb{P}} \mathbb{E}[X].$$

2. (*Mean-square convergence*) For any  $X$  such that  $\mathbb{E}[|X|^2] < \infty$ ,

$$\mathbb{E} \left[ \left| \frac{1}{T} \sum_{t=1}^T X_t(\omega) - \mathbb{E}[X] \right|^2 \right] \rightarrow 0.$$

- ▷ Once certain properties of  $\mathbb{P}$  and  $\mathbb{S}$  are satisfied, I can use any measurement function and it will be stationary.
- ▷ Why do we use  $\mathbb{S}$  instead of  $\mathbb{S}^{-1}$ ? Consider the example where  $\mathbb{S}(\omega_1) = \omega_2, \mathbb{S}(\omega_2) = \omega_1$ . Then

$$\begin{aligned} \mathbb{S}^{-1}(\{\omega_2\}) &= \{\omega_1, \omega_2\} \Rightarrow \{\omega_2\} \neq \mathbb{S}^{-1}(\{\omega_2\}) \\ \mathbb{S}(\{\omega_2\}) &= \{\omega_1\} \Rightarrow \{\omega_2\} \neq \mathbb{S}(\{\omega_2\}) \end{aligned}$$

so  $\{\omega_2\}$  is not invariant but using  $\mathbb{S}$  instead of  $\mathbb{S}^{-1}$  would seem like it satisfies the provided constraint.

### 2.3 Limiting Empirical Measures

- ▷ Limiting empirical measure: Given an event  $\Lambda \in \mathfrak{F}$  and a measure-preserving transformation  $\mathbb{S}$  for almost all  $\omega \in \Lambda_j$ , we define the limiting empirical measure  $Qr_j$  as

$$Qr_j(\Lambda)(\omega) := \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbf{1}_\Lambda(\mathbb{S}^t(\omega)) = \frac{P(\Lambda \cap \Lambda_j)}{P(\Lambda_j)}$$

Alternately, write

$$\mathbb{P}(\Lambda) = \sum_{j=0}^{\infty} \mathbb{P}(\Lambda \cap \Lambda_j) = \sum_{j=0}^{\infty} \left( \frac{P(\Lambda \cap \Lambda_j)}{P(\Lambda_j)} \right) \mathbb{P}(\Lambda_j) = \sum_{j=0}^{\infty} Qr_j(\Lambda) \mathbb{P}(\Lambda_j)$$

$Qr_j$  can be thought of as ergodic building blocks. They obey the law of large numbers in the usual sense.

- ▷  $Qr_j(\Lambda_j) = 1$  and  $Qr_j(\Lambda_k) = 0, j \neq k$ .
- ▷ If  $\mathbb{P}$  and  $\mathbb{S}$  is measure-preserving,  $(Qr_j, \mathbb{S})$  is measure-preserving and ergodic.

### 2.4 Ergodic Decomposition

- ▷ Start with  $Qr_j$  s and pick a prior distribution  $\pi_j$ s and construct

$$\mathbb{P}(\Lambda) = \sum_{j=0}^{\infty} Qr_j(\Lambda) \pi_j$$

and  $\mathbb{P}$  is measure-preserving.

## 2.5 Example

▷ Consider the VAR(1) of the form:

$$X_{t+1} = AX_t + BW_{t+1}, \quad W_{t+1} \sim_{iid} N(0, 1)$$

and  $A$  is stable i.e. eigenvalues have absolute values strictly less than 1.

\* *Mean*: Taking the expectation on both sides, the moments have to satisfy:

$$\mathbb{E}[X_{t+1}] = A\mathbb{E}[X_t]$$

and for stationarity, it must be that  $\mathbb{E}[X_t] = \mathbf{0}$ .

\* *Covariance*: Denoting it as  $\Sigma_t$ , we have

$$\Sigma_{t+1} = A\Sigma_t A' + BB'$$

Solving this using guess-and-verify:

$$\Sigma = \sum_{j=0}^{\infty} A^j BB' (A')^j$$

which is well-defined since  $A$  is stable.

\* *Stationarity*: We have stationarity:

$$X_t \sim N(0, \Sigma)$$

\* *Ergodic*: We have ergodicity here since we are conditioning on the fact that we know  $A$  and  $B$ .

▷ Consider:

$$X_t = \begin{bmatrix} X_t^1 \\ X_t^2 \end{bmatrix}, \quad A = \begin{bmatrix} A_{11} & A_{12} \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} B_1 \\ 0 \end{bmatrix}$$

so  $X_t^2$  is invariant for all  $t$  and suppose  $A_{11}$  is stable.

\* *Mean*: Taking the expectation, we have

$$\mathbb{E}[X_{t+1}^1] = A_{11}\mathbb{E}[X_t^1] + A_{12}\mathbb{E}[X_t^2] \Rightarrow (I - A_{11})\mathbb{E}[X_t^1] = A_{12}\mathbb{E}[X_t^2]$$

which yields

$$\mathbb{E}[X_t^1] = (I - A_{11})^{-1} A_{12}\mathbb{E}[X_t^2]$$

\* *Takeaway*:  $X_t^2$  is invariant, so depending on how you specify  $X_t^2$ , you get a different behavior for the mean of the process.  $X_t^2$  is the invariant event; so conditioning on  $X_t^2$  implies choosing different invariant sets.

\* It's only ergodic only if I announce that  $X_t^2$  is some specified number. The ergodic building block is precisely this specification of  $X_t^2$  that allows us to achieve ergodicity.



### 3 Markov Processes

#### 3.1 Finite-state Markov Chain

▷ Denote  $\mathbb{P}$  as a the transition matrix and  $p_{ij}$  as the transition probability and the state  $i$  is represented by a vector of zeros with only one at position  $i$ ,  $u_i$ . Denote  $X_t$  as the realization coordinate vectors.

▷ Consider a function of  $X_t$ ,  $f_{n \times 1}$ :

$$\mathbb{E}[f(X_{t+1}) | X_t = u] = u' \mathbb{P} f$$

or

$$\mathbb{E}[f(X_{t+1}) | X_t] = \mathbb{P} f$$

▷ Suppose we can find a solution to  $\mathbb{P} f = f$  i.e.

$$\mathbb{E}[f(X_{t+1}) | X_t] = f(X_t)$$

then  $\{X_t\}$  is stationary.

\* To see this, it suffices to show that given

$$f(X_{t+1}) = \underbrace{\mathbb{E}[f(X_{t+1}) | X_t]}_{=f(X_t)} + \epsilon_{t+1}$$

the variance of  $\epsilon_{t+1}$  is equal to zero i.e.  $\mathbb{E}[\epsilon_{t+1}^2] = 0$

\* Taking the expectation in the following manner:

$$\begin{aligned} & \mathbb{E}[(f(X_{t+1}) - \mathbb{E}[f(X_{t+1}) | X_t])^2] \\ &= \mathbb{E}[f(X_{t+1})^2] - 2\mathbb{E}[f(X_{t+1}) \mathbb{E}[f(X_{t+1}) | X_t]] + \mathbb{E}[\mathbb{E}[f(X_{t+1}) | X_t]^2] \\ &= \mathbb{E}[f(X_{t+1})^2] - 2\mathbb{E}[f(X_{t+1}) f(X_t)] + \mathbb{E}[f(X_t)^2] \\ &= \mathbb{E}[f(X_{t+1})^2] - 2\mathbb{E}[\mathbb{E}[f(X_{t+1}) f(X_t) | X_t]] + \mathbb{E}[f(X_t)^2] \\ &= \mathbb{E}[f(X_{t+1})^2] - 2\mathbb{E}\left[f(X_t) \underbrace{\mathbb{E}[f(X_{t+1}) | X_t]}_{:=f(X_t)}\right] + \mathbb{E}[f(X_t)^2] \\ &= \mathbf{0} \end{aligned}$$

## 4 Additive Functionals

### 4.1 Setup

- ▷ Martingale:  $\{Y_t\}$  is a martingale if  $\mathbb{E}[Y_{t+1}|\mathcal{F}_t] = Y_t$ . If  $Y_{t+1} - Y_t := X_{t+1}$ , then  $\mathbb{E}[X_{t+1}|\mathcal{F}_t] = 0$ .
- ▷ Additive Functional: A process  $\{Y_t\}$  is said to be an *additive functional* if it can be represented as

$$Y_{t+1} - Y_t = \kappa(X_t, W_{t+1})$$

for a (Borel measurable) function  $\kappa : \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}$ , or, equivalently,

$$Y_t = Y_0 + \sum_{j=1}^t \kappa(X_{j-1}, W_j),$$

where we initialise  $Y_0$  at some arbitrary (Borel measurable) function of  $X_0$ .

- ▷ Additive Martingale: An additive functional  $\{Y_t\}$  is said to be an *additive Martingale* if

$$\mathbb{E}[\kappa(X_t, W_{t+1}) | X_t] = 0.$$

\* A linear combination of additive functionals is an additive functional.

- ▷ Multiplicative Functional: Let  $\{Y_t\}$  be an additive functional that is described by  $Y_{t+1} - Y_t = \kappa(X_t, W_{t+1})$ . Then, we say that

$$\{M_t\} = \{\exp[Y_t]\}$$

is a *multiplicative functional* parameterised by  $\kappa$ .

- ▷ Multiplicative Martingale: A multiplicative functional  $\{M_t\}$  is said to be an *multiplicative Martingale* if

$$\mathbb{E}[M_{t+1} | \mathfrak{F}_t] = M_t.$$

### 4.2 Central Limit Theory

Billingsley (1961) proved the central limit theory result related to additive martingales without appealing to iid. Gordin (1969) then extended Billingsley's result to allow for temporally dependent increments.

**Theorem 4.1.** (Billingsley) Let  $\{Y_t\}_{t=0}^\infty$  be an additive martingale process whose increments  $Y_{t+1} - Y_t$  are stationary, ergodic, martingale differences:

$$\mathbb{E}[Y_{t+1} - Y_t | \mathfrak{F}_t] = 0.$$

Then,<sup>a</sup>

$$\frac{1}{\sqrt{t}} Y_t \xrightarrow{d} N\left(0, \mathbb{E}\left[(Y_{t+1} - Y_t)^2\right]\right).$$

**Corollary 4.1.** (Gordin) Suppose that  $\{Y_t\}$  is an additive functional, that  $\mathbb{T}^m$  is a strong contraction on  $\mathcal{N}$  for some  $m$ , and that  $\mathbb{E}[(\kappa(X_t, W_{t+1}))^2] < \infty$ . (basically the assumptions for martingale decomposition to work) and that  $\nu = 0$ . Then

$$\frac{1}{\sqrt{t}} Y_t \xrightarrow{d} N(0, \sigma^2),$$

where

$$\sigma^2 = \mathbb{E}\left[(\kappa_a(X_t, W_{t+1}))^2\right].$$

<sup>a</sup>Ergodicity can be dispensed with if we replace the variance by  $\mathbb{E}[(Y_{t+1} - Y_t)^2 | \mathfrak{J}]$ .

The variance formula

$$\sigma^2 = \lim_{t \rightarrow \infty} \frac{1}{t} \text{Var} [Y_t] = \mathbb{E} \left[ (\kappa_a (X_t, W_{t+1}))^2 \right]$$

shows how to take into account the temporal dependence of the increments  $(Y_{t+1} - Y_t)$  when computing the “long-run” volatility of the level  $Y_t$ . All that matters is the martingale component, and not the stationary  $g(X_t)$  component.

▷ Observation #1: In general,

$$\mathbb{E} [m_t^2] \neq \text{Var} [x_t]$$

▷ Observation #2: we have

$$\text{Var} \left[ \frac{1}{\sqrt{t}} \sum_{i=1}^t (x_i - \eta) \right] \rightarrow \text{Var} [m_t]$$

\* In practice, how should we compute the LHS? If we have a time series, we only obtain one observation of the term inside the bracket.

\* Alternatively, to estimate  $\text{Var} [m_t]$ , we can compute the spectral density at frequency zero of the process  $\{X_t - \eta\}$

▷ Observation #3: Setting  $y_0 = 0$ , we have

$$Y_t = \sum_{i=1}^t X_i$$

then

$$\begin{aligned} \frac{1}{t} Y_t &\rightarrow \eta \\ \frac{1}{\sqrt{t}} (Y_t - t\eta) &\rightarrow N(0, \mathbb{E} [m_t^2]) \end{aligned}$$

or alternatively:

$$\sqrt{t} \left( \frac{1}{t} Y_t - \eta \right) \rightarrow N(0, \mathbb{E} [m_t^2])$$

#### 4.2.1 Cointegration

Consider

$$Y_t = r_1 Y_t^1 + r_2 Y_t^2$$

where  $Y_t^1 = \eta_1 \{m_t^1\}$  and  $Y_t^2 = \eta_2 \{m_t^2\}$ . Then

$$m = r_1 m_1 + r_2 m_2$$

▷  $Y^1$  and  $Y^2$  are cointegrated if there exists  $r_1, r_2$  different from zero such that

$$r_1 m_t^1 + r_2 m_t^2 = 0$$

and then  $[r_1, r_2]^T$  is the cointegrating vector.  $[1, -1]^T$  was the cointegrating vector in the permanent income example in-class.

▷ Thinking about co-integration is like thinking about balanced growth path.

### 4.3 Martingale Decomposition Galore

#### 4.3.1 Martingale decomposition: Simple Version

*Example:* Suppose you have  $Y_{t+1} - Y_t = X_{t+1}$  where  $\{X_t\}$  is stationary and ergodic. Then  $Y_t = Y_0 + \sum_{j=1}^t X_j$ .

▷ Denoting  $\nu = \mathbb{E}[X_t]$ , define  $\bar{X}_t = \sum_{j=0}^{\infty} \mathbb{E}[X_{t+j} - \nu | \mathcal{F}_t]$ . We want this to be well-defined. (not obvious yet)

\* Define

$$\bar{X}_t = \sum_{j=0}^{\infty} \mathbb{E}[X_{t+j} - \nu | \mathcal{F}_t], \quad \tilde{X}_t = \sum_{j=1}^{\infty} \mathbb{E}[X_{t+j} - \nu | \mathcal{F}_t]$$

▷ Observation 1:  $\bar{X}_t - \tilde{X}_t = X_t - \nu$

▷ Observation 2:  $\mathbb{E}[\bar{X}_{t+1} | \mathcal{F}_t] = \tilde{X}_t$ . This can be shown by law of iterated expectations:

$$\begin{aligned} \mathbb{E}[\bar{X}_{t+1} | \mathcal{F}_t] &= \mathbb{E}\left[\sum_{j=0}^{\infty} \mathbb{E}[X_{t+1+j} - \nu | \mathcal{F}_{t+1}] | \mathcal{F}_t\right] \\ &= \sum_{j=0}^{\infty} \mathbb{E}[\mathbb{E}[X_{t+1+j} - \nu | \mathcal{F}_{t+1}] | \mathcal{F}_t] \\ &= \sum_{j=1}^{\infty} \mathbb{E}[\mathbb{E}[X_{t+j} - \nu | \mathcal{F}_{t+1}] | \mathcal{F}_t] \end{aligned}$$

Since  $\mathcal{F}_t \subseteq \mathcal{F}_{t+1}$ , we have

$$= \sum_{j=1}^{\infty} \mathbb{E}[X_{t+j} - \nu | \mathcal{F}_t] := \tilde{X}_t$$

▷ Observation 3:

$$\begin{aligned} X_{t+1} &= X_{t+1} - \nu + \nu = \bar{X}_{t+1} - \tilde{X}_{t+1} + \nu \\ &= \bar{X}_{t+1} - \tilde{X}_t + (\tilde{X}_t - \tilde{X}_{t+1}) + \nu \\ &= M_{t+1} + (\tilde{X}_t - \tilde{X}_{t+1}) + \nu \end{aligned}$$

where  $\mathbb{E}[M_{t+1} | \mathcal{F}_t] = 0$  since

$$\begin{aligned} \mathbb{E}[M_{t+1} | \mathcal{F}_t] &= \mathbb{E}[\bar{X}_{t+1} - \tilde{X}_t | \mathcal{F}_t] \\ &= \mathbb{E}[\bar{X}_{t+1} | \mathcal{F}_t] - \mathbb{E}[\tilde{X}_t | \mathcal{F}_t] \\ &= \tilde{X}_t - \tilde{X}_t = 0 \end{aligned}$$

▷ Observation 4: Using the result from observation 3, we have

$$\begin{aligned}
 Y_t &= Y_0 + \sum_{j=1}^t X_j \\
 &= Y_0 + \sum_{j=1}^t \left\{ \nu + M_j + (\tilde{X}_{j-1} - \tilde{X}_j) \right\} \\
 &= Y_0 + \nu t + \sum_{j=1}^t M_j + \tilde{X}_0 - \tilde{X}_t
 \end{aligned}$$

where  $Y_0 + \tilde{X}_0$  is invariant;  $\nu t$  is the time trend;  $\sum_j M_j$  is a martingale (the whole thing); and  $\tilde{X}_t$  is stationary.

▷ Note that  $M_j$  s are uncorrelated, which allows us to write:

$$Var \left[ \sum_{j=1}^t M_j \right] = t \mathbb{E} [M_t^2]$$

which is growing linear in time. To see why it's uncorrelated, consider  $k > j$  then

$$\mathbb{E} [M_j M_k] = \mathbb{E} [\mathbb{E} [M_k M_j | \mathcal{F}_{k-1}]] = \mathbb{E} [M_j \mathbb{E} [M_k | \mathcal{F}_{k-1}]] = 0$$

#### 4.3.2 Martingale decomposition: Markov's Version

Denote  $\{X_t\}$  to be Markov process and  $Y_{t+1} - Y_t = \phi(X_t)$  where  $\phi(X_t)$  has finite second moment.

1.  $\mathbb{T}$  as a contraction:

▷ Since we can characterize a Markov process with a transition matrix, write

$$\mathbb{T}\phi(X) = \mathbb{E} [\phi(X_{t+1}) | X_t = X]$$

which maps a function into a function. Note that we can go back and forth between the transition probabilities and the conditional expectations.

▷ Consider writing

$$\phi(X_{t+1}) = \mathbb{E} [\phi(X_{t+1}) | X_t] + \epsilon_{t+1}$$

which, by the least-squares theory, implies

$$\mathbb{E} [\phi(X_{t+1})^2] \geq \mathbb{E} [\{\mathbb{T}\phi(X_t)\}^2]$$

and thus  $\|\mathbb{T}\phi\| \leq \|\phi\|$  which implies that  $\mathbb{T}$  is a “weak” contraction. (We can't make it a strong contraction because the conditional expectation can be a constant).

▷ Define  $\mathcal{Z}_t = \{\phi \in \mathcal{L}^2 | \mathbb{E} [\phi(X)] = \mathbf{0}\}$  and assume a restriction that  $\mathbb{T}$  is a strong contraction on  $\mathcal{Z}_t$  i.e.  $\|\mathbb{T}\phi\| \leq \lambda \|\phi\|$ ,  $\lambda \in (0, 1)$ . This is not always satisfied, so we have to assume this.

▷ Suppose  $\sum_{j=0}^{\infty} \mathbb{E} [\phi(X_{t+j}) | X_t] = 0$  is mean zero. Then:

$$\sum_{j=0}^{\infty} \mathbb{E} [\phi(X_{t+j}) | X_t = X] = \sum_{j=0}^{\infty} \mathbb{T}^j \phi(X)$$

since LIE gives us:

$$\begin{aligned}
 \mathbb{E} [\phi (X_{t+2}) | X_t] &= \mathbb{E} [\mathbb{E} [\phi (X_{t+2}) | X_{t+1}, X_t] | X_t] \quad (\because \{X_t\} \subseteq \{X_{t+1}, X_t\}) \\
 &= \mathbb{E} [\mathbb{E} [\phi (X_{t+2}) | X_{t+1}] | X_t] \quad (\because \text{markov property}) \\
 &= \mathbb{T} [\mathbb{E} [\phi (X_{t+2}) | X_{t+1}]] \\
 &= \mathbb{T} [\mathbb{T} \phi (X_{t+1})] (X) \\
 &= \mathbb{T}^2 \phi (X)
 \end{aligned}$$

and the infinite sum above is well-defined and a strong contraction.

▷ This allows us to write:

$$(I - \mathbb{T})^{-1} = \sum_{j=0}^{\infty} \mathbb{T}^j$$

### 4.3.3 Martingale decomposition: VAR Example & Permanent Shocks

Setup is  $Y_{t+1} - Y_t = DX_t + FW_{t+1} + \nu$  where  $X_{t+1} = AX_t + BW_{t+1}$ ,  $A$  is stable, and  $W_{t+1} \sim N(O, I)$ . Furthermore, let

$$H_{t+1} = DX_t + FW_{t+1}$$

Then

$$\bar{H}_t := \mathbb{E} \left[ \left\{ \sum_{i=0}^{\infty} H_{t+i} \right\} | \mathcal{F}_t \right] = \dots = DX_{t-1} + FW_T + D \left[ (I - A)^{-1} X_t \right]$$

Computing the error term:

$$\begin{aligned}
 \bar{H}_{t+1} - \mathbb{E} [\bar{H}_{t+1} | \mathcal{F}_t] &= DX_t + FW_{t+1} + D \left[ (I - A)^{-1} X_{t+1} \right] - DX_t - D \left[ (I - A)^{-1} AX_t \right] \\
 &= FW_{t+1} + D \left[ (I - A)^{-1} BW_{t+1} \right]
 \end{aligned}$$

Going back to the VAR setup:

$$[X_{t+1}]_{n \times 1} = A_{n \times n} [X_t]_{n \times 1} + B_{n \times k} [W_{t+1}]_{m \times 1}$$

where

$$B = \begin{bmatrix} [B_1]_{n \times m} \\ O \end{bmatrix}$$

We want to infer  $BB'$ . To do this, let  $Q$  such that  $QQ' = I$ . then

$$B_1 W_{t+1} = B_1 Q Q' W_{t+1} = \tilde{B}_1 \tilde{W}_{t+1} \Rightarrow \tilde{B}_1 = B_1 Q, \tilde{W}_{t+1} = Q W_{t+1} \sim N(O, I)$$

Blanchard-Quah decomposition uses a permanent shock (in their paper this is the supply-side) to identify the shock of interest.

### 4.3.4 The Algorithm

We summarise the algorithm for extracting martingales below.

**Algorithm 1.** (Extracting martingales from additive functionals) The additive functional of interest is:

$$Y_{t+1} - Y_t = \kappa(X_t, W_{t+1}) \Leftrightarrow Y_t = Y_0 + \sum_{j=1}^t \kappa(X_{j-1}, W_j).$$

▷ Step 1: Construct  $\kappa_2$ .

$$\begin{aligned} \kappa_2(X_t, W_{t+1}) &:= \kappa(X_t, W_{t+1}) - \mathbb{E}[\kappa(X_t, W_{t+1}) | X_t] \\ &\equiv \kappa(X_t, W_{t+1}) - \bar{f}(X_t) \end{aligned}$$

▷ Step 2: Construct  $f(x)$ .

$$\begin{aligned} f(x) &:= \mathbb{E}[\kappa(X_t, W_{t+1}) | X_t = x] - \mathbb{E}[\kappa(X_t, W_{t+1})] \\ &\equiv \bar{f}(x) - \nu. \end{aligned}$$

▷ Step 3: Construct  $\kappa_1$  (assuming that  $\mathbb{T}$  is  $\rho$ -mixing).

$$\begin{aligned} \kappa_1(X_t, W_{t+1}) &= g(X_{t+1}) - g(X_t) + f(X_t), \\ g(X_t) &:= (\mathbb{I} - \mathbb{T})^{-1} f(X_t) = \sum_{j=0}^{\infty} \mathbb{T}^j f(X_t). \end{aligned}$$

▷ Step 4: Decompose  $\kappa$ .

$$\kappa(X_t, W_{t+1}) = \kappa_1(X_t, W_{t+1}) + \kappa_2(X_t, W_{t+1}) + g(X_t) - g(X_{t+1}) + \nu.$$

▷ Step 5: Substitute into  $Y_t$ .

$$Y_t = \underbrace{Y_0 + g(X_0)}_{\text{Invariant}} + \underbrace{t\nu}_{\text{Time trend}} - \underbrace{g(X_t)}_{\text{Stationary}} + \underbrace{\sum_{j=1}^t \kappa_a(X_{j-1}, W_j)}_{\text{Martingale}}.$$

#### 4.4 TA Session: Permanent vs Transient Shocks

We illustrate them using the VAR example:

$$\begin{aligned} X_{t+1} &= AX_t + BW_{t+1} \\ Y_{t+1} - Y_t &= \nu + D^T X_t + F^T W_{t+1} \end{aligned}$$

where  $A$  is stable and  $X_0 = 0$ .  $X_{t+1}$  is a  $k \times 1$  vector;  $Y_t$  is a scalar.

▷ In class, we saw the following decomposition:

$$Y_t = Y_0 + \nu t + \underbrace{\sum_{j=1}^t \left( F^T + D^T (I - A)^{-1} B \right) W_j - D^T (I - A)^{-1} X_t}_{= \text{martingale component}}$$

▷ Using the lag operator, rewrite the process as

$$(I - AL) X_t = BW_t$$

and plug this back into the decomposition:

$$Y_t = Y_0 + \nu t + \sum_{j=1}^t \left( F^T + D^T (I - A)^{-1} B \right) W_j - D^T (I - A)^{-1} \underbrace{(I - AL)^{-1} BW_t}_{(1)}$$

Note that (1) is

$$(I + AL + A^2L^2 + \dots) BW_t = BW_t + ABW_{t-1} + A^2BW_{t-2} + \dots$$

so we have

$$Y_t = Y_0 + \nu t + \sum_{j=1}^t \left( F^T + D^T (I - A)^{-1} B \right) W_j - D^T (I - A)^{-1} \sum_{j=0}^{t-1} A^j BW_{t-j}$$

▷ Assume  $W_1 = 0$  and  $W_t = 0, \forall t \geq 2$ . Then

$$Y_t = Y_0 + \nu t + \underbrace{\left( F^T + D^T (I - A)^{-1} B \right) W_1}_{\text{permanent part}} - \underbrace{D^T (I - A)^{-1} A^{t-1} BW_1}_{\text{transient part}}$$

The transiency comes from the decay that comes from  $A^{t-1}$ .

- \* Permanent part is just a linear combination of  $W_1$ .
- \* If  $W_1$  is *orthogonal* to  $F^T + D^T (I - A)^{-1} B$ , we call it the *transitory* shock.
- \* If  $W_1$  is *parallel* to  $F^T + D^T (I - A)^{-1} B$ , we call it the *permanent* shock.

- ▷ For a transitory impulse response, the impulse should just decay exponentially.
- ▷ For all shocks that are not transitory, the impulse response will converge to
- ▷ The permanent shock can be found by fixing the magnitude of the impulse response and finding the  $W_1$  that gives the largest steady-state.
- \* Even if the shock is entirely parallel, you still have a convergence over time to the steady state.

## 4.5 TA Session: Small Shock Approximation a la Lombardo and Uhlig (2018)

### 4.5.1 Setup

Consider the following setup:

$$\begin{aligned} X_{t+1} &= AX_t + BW_{t+1} \\ y_{t+1} - y_t &= \nu + D^T X_t + F^T W_{t+1} \end{aligned}$$

where  $y_t := \log Y_t$ . We will illustrate how the previous technique can be applied in a macroeconomic framework.

- ▷ Technology: transfer 1 unit time  $t$  good to  $\exp(\rho)$  units of time  $t + 1$  good. The  $\{Y_t\}$ s are the fruits.



▷ The feasibility constraint:

$$K_{t+1} + C_t = \exp(\rho) K_t + Y_t$$

This is equivalent to Prof. Stokey's formulation:

$$K_{t+1} + C_t = A_t K_t + (1 - \delta) K_t$$

▷ The RA maximizes

$$\mathbb{E} \left[ \sum_{j=0}^{\infty} \exp(-\delta t) \log C_t \right]$$

▷ The Euler Equation:

$$U'(C_t) = \mathbb{E} [\exp(\rho) \exp(-\delta) u'(C_{t+1}) | X_t]$$

which yields

$$1 = \mathbb{E} \left[ \exp(-\delta + \rho) \frac{C_t}{C_{t+1}} | X_t \right]$$

#### 4.5.2 Solving for Response in Consumption

We want to linearize the  $\mathbb{E}[\cdot]$ .

▷ To solve this, define

$$\hat{K}_t = \frac{K_t}{Y_t}, \quad \hat{C}_t = \log C_t - \log Y_t$$

and re-write the feasibility constraint:

$$\hat{K}_{t+1} \exp(\log Y_{t+1} - \log Y_t) + \exp(\hat{C}_t) - \exp(\rho) \hat{K}_t - 1 = 0$$

and the Euler equation:

$$\exp(-\delta + \rho) \mathbb{E} \left[ \exp \left( - \left( \hat{C}_{t+1} - \hat{C}_t \right) - (\log Y_{t+1} - \log Y_t) \right) | X_t \right] - 1 = 0$$

Essentially, we are detrending to get the stationary distribution.

▷ Now consider perturbing the system by  $q$  i.e. changing the exposure of  $Y_t$ s to the stochastic component:

$$\begin{aligned} X_{t+1} &= AX_t + BW_{t+1} \\ y_{t+1}(q) - y_t(q) &= \nu + [D^T X_t + F^T W_{t+1}] q \end{aligned}$$

This allows us to reformulate the previous variables as a function of  $q$

$$[1] : \hat{K}_{t+1}(q) \exp(\log Y_{t+1}(q) - \log Y_t(q)) + \exp(\hat{C}_t(q)) - \exp(\rho) \hat{K}_t(q) - 1 = 0$$

$$[2] : \exp(-\delta + \rho) \mathbb{E} \left[ \exp \left( - \left( \hat{C}_{t+1}(q) - \hat{C}_t(q) \right) - (\log Y_{t+1}(q) - \log Y_t(q)) \right) | X_t \right] - 1 = 0$$

▷ Consider a Taylor expansion around 0:

$$\begin{aligned} \hat{C}_t(q) &\approx \hat{C}_t(0) + \hat{C}'_t(q) \\ \hat{K}_{t+1}(q) &\approx \hat{K}_{t+1}(0) + \hat{K}'_{t+1}(q) \end{aligned}$$

The reason we do it around 0 is because  $\hat{C}_t(0)$  results in a deterministic fruits process  $\{Y_t\}$  and thus it is very easy to compute.

▷ Note that the individual components above are processes, not numbers. Similarly to the macro class, we want to do

$$\hat{C}_t = C \left( X_t, \hat{K}_t \right), \quad \hat{K}_{t+1} = K \left( X_t, \hat{K}_t \right)$$

where  $X_t$  is the exogenous state and  $K_t$  is the endogenous state.

Now define a new function

$$F_1 \left( \hat{K}_{t+1}(q), \hat{C}_t(q), \hat{K}_t(q), \Delta y_{t+1}(q) \right) \equiv [1]$$

where  $\Delta y_{t+1}(q) = y_{t+1}(q) - y_t(q)$

▷ Since  $F_1$  is equal to zero for all  $q$ , we have

$$F_1(q) \approx F_1|_{q=0} + q \frac{\partial F_1}{\partial q} \Big|_{q=0} \approx 0$$

as well. Similar argument holds for  $F_2 = 0$ :

$$F_2(q) \approx F_2|_{q=0} + q \frac{\partial F_2}{\partial q} \Big|_{q=0} \approx 0$$

▷ Obtaining  $F_1|_{q=0}$ : Making the assumption that  $\delta = \rho - \nu$ , the economy has a steady state of:

$$\hat{C}_t(0) = 0, \quad \hat{K}_{t+1}(0) = 0$$

which yields

$$\begin{aligned} \log C_t - \log Y_t &= 0 \Rightarrow C_t = Y_t \\ \hat{K}_{t+1}(0) &= \frac{K_{t+1}}{Y_{t+1}} = 0 \end{aligned}$$

so you consume fruit everyday and save nothing.

▷ Obtaining  $\partial F_1 / \partial q$ :

$$\frac{\partial F_1}{\partial q} \Big|_{q=0} = \hat{C}'_t(0) \frac{\partial F_1}{\partial \hat{C}_t} \Big|_{q=0} + \hat{K}'_{t+1}(0) \frac{\partial F_1}{\partial \hat{K}_{t+1}} \Big|_{q=0} + \hat{K}'_t(0) \frac{\partial F_1}{\partial \hat{K}_t} \Big|_{q=0} + \Delta y'_{t+1}(0) \frac{\partial F_1}{\partial \Delta y_{t+1}} \Big|_{q=0} = 0$$

\* Note that we already know

$$\frac{\partial F_1}{\partial \hat{C}_t} \left( \hat{K}_{t+1}(q), \hat{C}_t(q), \hat{K}_t(q), \Delta y_{t+1}(q) \right) \Big|_{q=0}$$

since we've computed the relevant quantities evaluated at zero in the previous step. The similar argument follows for the other derivatives.

\* Deriving analogously for  $F_2$ , we have two linear equations of  $\hat{K}'_{t+1}(0)$  and  $\hat{C}'_t(0)$ .

\* The term with  $\Delta y'_{t+1}$  is good since

$$\frac{\partial F_1}{\partial \Delta y_{t+1}} \Big|_{q=0} = 0$$

and the term vanishes.

\* The term with  $\hat{K}'_t(0)$  is also good since the derivative is simply  $\exp(\rho)$ .

Going through a similar process with  $F_2$ , we obtain:

$$[3] : \hat{K}_{t+1}^1 \exp(\nu) + \hat{C}_t^1 - \exp(\rho) \hat{K}_t' = 0$$

$$[4] : \mathbb{E} \left[ \left( \hat{C}_{t+1}^1 - \hat{C}_t^1 \right) + \Delta y_{t+1}' \right] = 0$$

To solve this, guess and verify:

$$\hat{C}_t^1 := C_t'(q=0) = M X_t + \Gamma_K \hat{K}_t^1$$

- ▷ Plug this into [3] and express  $\hat{K}_{t+1}$  as a function of  $\hat{K}_t$  and  $\hat{X}_t$ .
- ▷ Plug this into [4] and replace  $X_{t+1}$  as a function of  $X_t$  and solve for  $M$  and  $\Gamma_K$ .

The resulting solution is

$$\begin{aligned} \hat{C}_t' &= \lambda D' (I - \lambda A)^{-1} X_t + \{\exp(\rho)\} \\ \hat{K}_{t+1}' &= \hat{K}_t' - \exp(-\nu) \lambda D' (I - \lambda A)^{-1} X_t \end{aligned}$$

where  $\lambda = \exp(\nu - \rho)$ . We also obtain

$$\hat{C}_{t+1}^1 - \hat{C}_t^1 = -D^T X_t + \lambda D^T (I - \lambda A)^{-1} B W_{t+1}$$

This allows us to compute the log consumption progress:

$$\begin{aligned} \log C_{t+1} - \log C_t &= -D^T X_t + \lambda D^T (I - \lambda A)^{-1} B W_{t+1} + \nu + D^T X_t + F^T W_{t+1} \\ &= \left( \lambda D^T (I - \lambda A)^{-1} B + F^T \right) W_{t+1} + \nu \end{aligned}$$

This is the function that Professor Hansen plotted in class. He provided two plots – one is the permanent shock to the log  $Y_t$  process and one is the transitory shock to the log  $Y_t$  process.

- ▷ If  $\lambda$  is really close to one, then the permanent shock which is parallel to

$$D^T (I - \lambda A)^{-1} B + F^T$$

will be almost parallel to

$$\lambda D^T (I - \lambda A)^{-1} B + F^T$$

in which case the permanent shock will have a big impact on consumption.

- ▷ There is no  $D$  which is why you see the impulse response as a constant. This is contrast with the general response with the long-term convergence (adjustment takes time due to the transitory part, which was  $D[\cdot \cdot \cdot]$ ). If  $D = 0$ , then it will be just a straight line.

### 4.5.3 Comparison with Log-linearization

If you want to attain the second order in log-linearization, you have to deal with

$$\hat{X}_{t+1} = a \hat{X}_t + b \hat{X}_t^2 + \dots$$

which fucks with stationarity. But for this methodology, even if you go to second order:

$$\hat{X}_{t+1} = \hat{X}_{t+1}(0) + [\dots] \hat{X}_t'(0) + [\dots] \hat{X}_t''(0)$$

you are taking the square of the derivative so it's okay.

## 5 Likelihood Processes (Ch. 8)

### 5.1 Likelihood Constructions

Let  $\{W_t\}$  be a process of shocks satisfying  $\mathbb{E}[W_{t+1}|\mathcal{F}_t] = 0$ . Let  $\{X_t\}$  be a Markov process such that  $X_{t+1} = \phi(X_t, W_{t+1})$ . We also have data that we observe, denoted as  $Y_t$ , such that

$$Y_{t+1} - Y_t = \kappa(X_t, W_{t+1})$$

We make the following assumptions:

- ▷ There exists a function  $\chi$  such that  $W_{t+1} = \chi(X_t, Y_{t+1} - Y_t)$
- ▷  $Y_{t+1} - Y_t$  has a density  $\psi(\cdot|x)$  with respect to measure  $\tau$  conditioned on  $X_t = x$ .

The example below illustrates how to construct  $\chi$  and  $\psi$ .

**Example 5.1.** (8.1.4 of Lecture Notes) Suppose  $X_{t+1} = AX_t + BW_{t+1}$  and  $Y_{t+1} - Y_t = DX_t + FW_{t+1}$  where  $A$  is a stable matrix,  $\{W_{t+1}\}_{t=0}^\infty$  is an i.i.d. sequence with mean 0 and covariance  $I$ . Also assume  $F$  is non-singular. Then applying  $F^{-1}$  to both sides:

$$W_{t+1} = F^{-1}(Y_{t+1} - Y_t) - F^{-1}DX_t := \chi(X_t, Y_{t+1} - Y_t)$$

And the density is then

$$Y_{t+1} - Y_t \sim N(DX_t, FF') := \psi(\cdot|x)$$

The likelihood function conditioned on  $X_0$ :

$$\begin{aligned} L_t &= \prod_{j=1}^t \psi(Y_j - Y_{j-1}|X_{j-1}) \\ \Rightarrow \ell_t &= \sum_{j=1}^t \log(\psi(Y_j - Y_{j-1}|X_{j-1})) \end{aligned}$$

and  $\log L_t := \ell_t$  are stationary increments. Both of these are stochastic processes.

### 5.2 Likelihood Ratios

We show that the likelihood ratio process is a positive martingale.

1. Construct  $L_t(\theta)$  and  $L_t(\theta_0)$  recursively as above
2. Decompose the likelihood in the following manner:

$$\frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} = \frac{L_t(\theta)}{L_t(\theta_0)} \cdot \left[ \frac{\psi(Y_{t+1} - Y_t|X_t, \theta)}{\psi(Y_{t+1} - Y_t|X_t, \theta_0)} \right]$$

3. Taking the expectation of the LHS:

$$\begin{aligned} \mathbb{E} \left[ \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} | X_t = x, \theta_0 \right] &= \frac{L_t(\theta)}{L_t(\theta_0)} \mathbb{E} \left[ \frac{\psi(Y_{t+1} - Y_t|X_t, \theta)}{\psi(Y_{t+1} - Y_t|X_t, \theta_0)} | X_t = x, \theta_0 \right] \\ &= \frac{L_t(\theta)}{L_t(\theta_0)} \int_Y \frac{\psi(y^*|x, \theta)}{\psi(y^*|x, \theta_0)} \psi(y^*|x, \theta_0) \tau(dy) \\ &= \frac{L_t(\theta)}{L_t(\theta_0)} \int_Y \psi(y^*|x, \theta) \tau(dy) \\ &= \frac{L_t(\theta)}{L_t(\theta_0)} \end{aligned}$$

and thus we verify that the likelihood ratio process is a multiplicative (positive) martingale under the  $\theta_0$  probability model.

Next we show that the log likelihood ratio process is a super-martingale. This is easy since from Jensen's inequality, we have

$$\mathbb{E} \left[ \log \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} | X_t = x, \theta_0 \right] \leq \log \left[ \frac{L_t(\theta)}{L_t(\theta_0)} \right]$$

Now using the above result, we have define  $v(\theta)$  as

$$\begin{aligned} v(\theta) &:= \mathbb{E} [\log \psi(Y_{t+1} - Y_t | X_t, \theta) - \log \psi(Y_{t+1} - Y_t | X_t, \theta_0) | \theta_0] \\ &= \mathbb{E} \left[ \log \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} - \log \frac{L_t(\theta)}{L_t(\theta_0)} | \theta_0 \right] \\ &= \mathbb{E} \left[ \underbrace{\mathbb{E} \left[ \log \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} - \log \frac{L_t(\theta)}{L_t(\theta_0)} | X_t = x, \theta_0 \right]}_{\leq 0} | \theta_0 \right] \leq 0 \\ \Rightarrow v(\theta) &= \mathbb{E} [\log \psi(Y_{t+1} - Y_t | X_t, \theta) - \log \psi(Y_{t+1} - Y_t | X_t, \theta_0) | \theta_0] \leq 0 \end{aligned}$$

since the conditional expectation is  $\leq 0$ . Typically, it is  $< 0$ . This leads to

$$\frac{1}{t} [\log L_t(\theta) - \log L_t(\theta_0)] \rightarrow v(\theta) < 0$$

So the log likelihood process will tend to infinity, implying that the ratio of likelihood tends to zero.

▷ First, we can show that  $v(\theta_0) \geq v(\theta), \forall \theta \in \Theta$ .

\* Starting with

$$\mathbb{E} \left[ \log \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} | X_t = x, \theta_0 \right] \leq \log \left[ \frac{L_t(\theta)}{L_t(\theta_0)} \right]$$

\* Then:

$$\begin{aligned} 0 &\geq \mathbb{E}_{\theta_0} [(\log L_t(\theta) - \log L_{t-1}(\theta)) - (\log L_t(\theta_0) - \log L_{t-1}(\theta_0)) | \mathfrak{F}_{t-1}] \\ &= \mathbb{E}_{\theta_0} [\log L_t(\theta) - \log L_{t-1}(\theta) | \mathfrak{F}_{t-1}] - \mathbb{E} [\log L_t(\theta_0) - \log L_{t-1}(\theta_0) | \mathfrak{F}_{t-1}] \\ &= \mathbb{E}_{\theta_0} [g(X_t | \theta) - g(X_{t-1} | \theta) + \nu(\theta) + \kappa_a(X_{t-1}, W_t | \theta) | \mathfrak{F}_{t-1}] \\ &\quad - \mathbb{E}_{\theta_0} [g(X_t | \theta_0) - g(X_{t-1} | \theta_0) + \nu(\theta_0) + \kappa_a(X_{t-1}, W_t | \theta_0) | \mathfrak{F}_{t-1}] \\ [*] &= \nu(\theta) - \nu(\theta_0) + \mathbb{E}_{\theta_0} [g(X_t | \theta) - g(X_t | \theta_0) | \mathfrak{F}_{t-1}] \\ \Rightarrow 0 &\geq \nu(\theta) - \nu(\theta_0), \forall \theta \in \Theta, \end{aligned}$$

where in step  $[*]$ , we use the fact that conditional on  $\mathfrak{F}_{t-1}$ ,  $g(X_{t-1} | \theta) = g(X_{t-1} | \theta_0)$ , and in the last step, we take unconditional expectation of both sides and use the fact that the unconditional expectation of  $g$ 's are zero. Finally, recall that the original inequality were strict unless the two log-likelihoods are equal; i.e. inequality above holds with equality if and only if  $\theta = \theta_0$  (with probability one).

▷ Second, the above result naturally implies that the solution to

$$\max_{\theta \in \Theta} v(\theta)$$

is  $\theta = \theta_0$ .

▷ Third, the solution to the above problem also solves the problem below:

$$\max_{\theta \in \Theta} \hat{v}_t(\theta) = \max_{\theta \in \Theta} \lim_{t \rightarrow \infty} \frac{1}{t} \log L_t(\theta)$$

### 5.3 Score Processes

Suppose you want to solve the following problem conditionally, not unconditionally:

$$\max_{\theta} \int_y \log(\psi(y^*|x, \theta)) \psi(y^*|x, \theta_0) \tau(dy^*) := \mathbb{E}_{\theta_0} [\log(\psi(y^*|x, \theta))]$$

▷ Since we had this inequality from before::

$$\mathbb{E} \left[ \log \frac{L_{t+1}(\theta)}{L_{t+1}(\theta_0)} | X_t = x, \theta_0 \right] \leq \log \left[ \frac{L_t(\theta)}{L_t(\theta_0)} \right]$$

then the parameter  $\theta_0$  necessarily solves the above maximization problem since

$$\mathbb{E}_{\theta_0} [\log(\psi(y^*|x, \theta))] \leq \log(\mathbb{E}_{\theta_0} [\psi(y^*|x, \theta)])$$

▷ Suppose that we can differentiate under the integral sign and get to the first-order condition (for interior solutions)

$$\begin{aligned} 0 &= \int_{\mathcal{Y}} \left[ \frac{d}{d\theta} \log \psi(y^*|x, \theta_o) \right] \psi(y^*|x, \theta_o) \tau(dy^*) \\ &= \mathbb{E}_{\theta_o} \left[ \frac{d}{d\theta} \log \psi(y^*|x, \theta_o) | X_t, \theta_o \right]. \end{aligned}$$

Define the score process  $\{S_t\}$  as

$$S_t = \frac{d}{d\theta} \log L_t(\theta) \Big|_{\theta=\theta_o} = \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) + \frac{d}{d\theta} \log L_0(\theta_o).$$

▷ Using this definition, note that

$$S_{t+1} - S_t = \frac{d}{d\theta} \log \psi(Y_{t+1} - Y_t | X_t, \theta_o).$$

thus implying

$$\mathbb{E}_{\theta_o} [S_{t+1} - S_t | X_t, \theta_o] = 0.$$

Thus,  $\{S_t\}$  is a martingale.

▷ Property #1:

$$\frac{1}{\sqrt{t}} S_t \xrightarrow{d} N(0, \mathbf{V}),$$

where  $\mathbf{V}$  is Fisher information matrix

$$\mathbf{V} = \mathbb{E} \left[ (S_{t+1} - S_t) (S_{t+1} - S_t)^T | \theta_0 \right] = \mathbb{E}_{\theta_o} \left[ \left( \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) \right)^2 \right].$$

\* To see this, it helps to write

$$\begin{aligned} S_t &= (S_t - S_{t-1}) + (S_{t-1} - S_{t-2}) + \cdots + (S_1 - S_0) + S_0 \\ &= S_0 + \sum_{j=1}^t (S_j - S_{j-1}) = S_0 + \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) \\ \Rightarrow \frac{1}{\sqrt{t}} S_t &= \frac{1}{\sqrt{t}} S_0 + \sqrt{t} \frac{1}{t} \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o). \end{aligned}$$

▷ Property #2: With some additional regularity conditions, the following result is typical:<sup>2</sup>

$$\sqrt{t} (\theta_t^{MLE} - \theta_o) \xrightarrow{d} N(0, \mathbf{V}^{-1}),$$

where  $\theta_t^{MLE}$  maximises the log-likelihood function  $\log L_t(\theta)$ .

## 5.4 Nuisance Parameters

Suppose there is a vector  $\theta$  of unknown parameters but we are interested in information about only one of the components of the parameter vector. Call the first component of  $\theta$  the “parameter of interest”  $\theta_o$ , and other components the “nuisance parameters”  $\vartheta_0 := \tilde{\theta}_0$ .

1. Write the multivariate score process as

$$\{\mathbf{S}_{t+1}\}_{t=0}^{\infty} = \left\{ \begin{bmatrix} S_{t+1} \\ \tilde{\mathbf{S}}_{t+1} \end{bmatrix} \right\}_{t=0}^{\infty},$$

where  $\{S_{t+1}\}$  is the partial derivative of the log-likelihood with respect to  $\theta$  and  $\{\tilde{\mathbf{S}}_{t+1}\}$  is the partial derivative with respect to  $\vartheta_o$ .

2. Run a population regression

$$\begin{aligned} S_{t+1} - S_t &= \beta' (\tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t) + U_{t+1} \\ \Rightarrow \begin{bmatrix} S_{t+1} - S_t \\ \tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t \end{bmatrix} &= \begin{bmatrix} 1 & \beta' \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} U_{t+1} \\ \tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t \end{bmatrix}. \end{aligned}$$

---

<sup>2</sup>From Bonhomme's class.  $\hat{\theta}_t$  satisfies the first-order condition:

$$\frac{1}{\sqrt{t}} S_t = \sqrt{t} \frac{1}{t} \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \hat{\theta}_t) = 0.$$

Taylor expansion around  $\theta_o$  yields

$$0 \approx \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) + \sum_{j=1}^t \frac{d}{d\theta d\theta'} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \tilde{\theta}_t) (\hat{\theta}_t - \theta_o),$$

where  $\tilde{\theta}_t$  lies in between  $\hat{\theta}_t$  and  $\theta_o$  (component by component). Rearranging yields

$$\begin{aligned} \hat{\theta}_t - \theta_o &\approx - \left[ \sum_{j=1}^t \frac{d}{d\theta d\theta'} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \tilde{\theta}_t) \right]^{-1} \left( \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) \right) \\ \Rightarrow \sqrt{t} (\hat{\theta}_t - \theta_o) &\approx - \left[ \frac{1}{t} \sum_{j=1}^t \frac{d}{d\theta d\theta'} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \tilde{\theta}_t) \right]^{-1} \left( \sqrt{t} \frac{1}{t} \sum_{j=1}^t \frac{d}{d\theta} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) \right) \\ &= - \left[ \frac{1}{t} \sum_{j=1}^t \frac{d}{d\theta d\theta'} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \tilde{\theta}_t) \right]^{-1} \frac{1}{\sqrt{t}} S_t. \end{aligned}$$

The information matrix identity (see end of this section) gives us that

$$-\mathbb{E}_{\theta_o} \left[ \frac{d}{d\theta d\theta'} \log \psi(Y_j - Y_{j-1} | X_{j-1}, \theta_o) \right] = \mathbf{V}.$$

So that, in fact

$$\sqrt{t} (\hat{\theta}_t - \theta_o) \xrightarrow{d} N(0, \mathbf{V}^{-1} \mathbf{V} \mathbf{V}^{-1}) \stackrel{d}{=} N(0, \mathbf{V}^{-1}).$$

## 3. Construct the matrix

$$\begin{aligned}\hat{\mathbf{V}} &= \mathbb{E} \left[ \begin{bmatrix} S_{t+1} - S_t \\ \tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t \end{bmatrix} \begin{bmatrix} S_{t+1} - S_t & \tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t \end{bmatrix} \right] \\ &= \underbrace{\begin{bmatrix} 1 & \beta' \\ \mathbf{0} & \mathbf{1} \end{bmatrix}}_{=\mathbf{A}} \underbrace{\begin{bmatrix} \mathbb{E}[U_{t+1}^2] & \mathbf{0} \\ \mathbf{0} & \mathbb{E}[(\tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t)(\tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t)'] \end{bmatrix}}_{=\mathbf{B}} \underbrace{\begin{bmatrix} 1 & \mathbf{0} \\ \beta & \mathbf{1} \end{bmatrix}}_{=\mathbf{A}'}\end{aligned}$$

then

$$\hat{\mathbf{V}}^{-1} = \begin{bmatrix} 1 & -\beta' \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \frac{1}{\mathbb{E}[U_{t+1}^2]} & \mathbf{0} \\ \mathbf{0} & \left( \mathbb{E}[(\tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t)(\tilde{\mathbf{S}}_{t+1} - \tilde{\mathbf{S}}_t)'] \right)^{-1} \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0} \\ -\beta & \mathbf{1} \end{bmatrix}$$

and the (1,1) element, which is our main point of interest is given by

$$\hat{v}_{11}^{-1} = \frac{1}{\mathbb{E}_{\theta_0}[U_{t+1}^2]}.$$

which is referred to as the Fisher Information for estimating  $\theta_0$  with  $\tilde{\theta}_0$  unknowns.

## 4. Notice that

$$\mathbb{E}[(S_{t+1} - S_t)^2] = \text{Var}[S_{t+1} - S_t] \geq \text{Var}[U_{t+1}] = \mathbb{E}[U_{t+1}^2] = \hat{v}_{11},$$

where  $\hat{v}_{11}$  is information for  $\theta_o$

- ▷ Since the left-hand side is the information for  $\theta_o$  when we know the nuisance parameters, we see that the likelihood function contains more information about  $\theta$  when  $\vartheta$  is known to be  $\vartheta_o$  than when  $\theta$  and  $\vartheta$  are unknown.

The punchline is that the more nuisance parameters you have, the more risk of losing information. We only get a nice CLT for additive martingales but not for multiplicative martingales.



## 6 Learning (Ch. 9)

### 6.1 Learning about discrete states

Suppose that  $\{\mathbf{X}_t\}$  evolves as an  $n$ -state Markov process with (a known) transition matrix  $\Pi$ . Crucially, we assume that  $\mathbf{X}_t$  is unknown. However, there is a vector of signals denoted by  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  with density  $\psi_i(\mathbf{y}^*)$  if the state  $i$  is realized (i.e. if  $\mathbf{X}_t$  is the  $i$ th coordinate vector). We want to compute the probability that the state is  $i$  given the signal history; i.e. we want to compute the vector of conditional probabilities:

$$\mathbf{Q}_t = \begin{bmatrix} q_t^1 \\ q_t^2 \\ \vdots \\ q_t^n \end{bmatrix} = \mathbb{E} [\mathbf{X}_t | \mathbf{Y}^t, \mathbf{Q}_0] = \mathbb{P} [\mathbf{X}_t | \mathbf{Y}^t, \mathbf{Q}_0],$$

where  $\mathbf{Q}_0$  is the vector of initial probabilities and  $q_t^i$  gives the probability that the current state is the  $i$ th state. We will do so recursively; i.e. we will derive the expression for

$$\mathbf{Q}_{t+1} = \mathbb{P} (\mathbf{X}_{t+1} | \mathbf{Q}_t, \mathbf{Y}^{t+1}, \mathbf{Q}_0).$$

Using Bayes' rule,<sup>3</sup> we can write  $\mathbf{Q}_{t+1}$  as

$$\begin{aligned} \mathbf{Q}_{t+1} &= \mathbb{P} (\mathbf{X}_{t+1} | \mathbf{Q}_t, \mathbf{Y}_{t+1} - \mathbf{Y}_t, \mathbf{Y}^t, \mathbf{Q}_0) \\ &= \frac{\mathbb{P} (\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t, \mathbf{Y}^t, \mathbf{Q}_0)}{\mathbb{P} (\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t, \mathbf{Y}^t, \mathbf{Q}_0)} \\ &= \frac{\mathbb{P} (\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)}{\mathbb{P} (\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)}, \end{aligned}$$

where  $\mathbf{Q}_0$  is given. Our goal is therefore to derive expressions for the numerator and the denominator of the expression above.

**Step 1: Find the joint distribution for  $(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t) | \mathbf{X}_t$ .** (Further past values of  $\mathbf{X}_t$  is not relevant due to the Markovian structure). We begin by assuming that the state is known. The probability that tomorrow's state is the  $j$ th coordinate vector is given by  $\pi_{ij}$  for all  $j = 1, 2, \dots, n$ ; i.e.  $(\pi_{i1}, \pi_{i2}, \dots, \pi_{in})' = \Pi' \mathbf{X}_t$ . So,

$$\mathbb{P} (\mathbf{X}_{t+1} | \mathbf{X}_t = \mathbf{x}_i) = \Pi' \mathbf{X}_t = \begin{bmatrix} \pi_{i1} \\ \pi_{i2} \\ \vdots \\ \pi_{in} \end{bmatrix}$$

gives the probability distribution over the next-period state.

Let  $\text{vec} [\psi_i(\mathbf{y}^*)]$  denote the column vector of  $\psi_i(\mathbf{y}^*)$ 's. Given  $\mathbf{X}_t$  as the  $i$ th coordinate vector, the density for  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  is given by

$$\psi_i(\mathbf{y}^*) = \mathbf{X}_t' \text{vec} [\psi_i(\mathbf{y}^*)].$$

Since  $\mathbf{X}_{t+1}$  and  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  are independent conditional on  $\mathbf{X}_t$  by assumption, the joint density conditioned on  $\mathbf{X}_t$  is given by the product of the two probabilities; i.e.

$$\begin{aligned} \mathbb{P} (\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{X}_t) &= \mathbb{P} (\mathbf{X}_{t+1} | \mathbf{X}_t) \mathbb{P} (\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{X}_t) \\ &= \Pi' \mathbf{X}_t \mathbf{X}_t' \text{vec} [\psi_i(\mathbf{y}^*)] \end{aligned}$$

<sup>3</sup>Recall that  $\mathbb{P}(A|B) = \mathbb{P}(A, B) / \mathbb{P}(B)$ . Here, let  $A = \mathbf{X}_{t+1}$  and  $B = \mathbf{Y}_{t+1} - \mathbf{Y}_t$  and condition everything on  $\mathbf{Q}_t$ .

**Step 2: Find the joint distribution of  $(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t) | \mathbf{Q}_t$ .** In step 1, we assumed that the state is known. But the premise here is that the states are hidden and so, in fact, we can only condition on  $\mathbf{Q}_t$ , which is the vector of probabilities of the current state. We obtain the joint distribution of  $(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t) | \mathbf{Q}_t$  by taking an “average” of (??) (conditioned on  $\mathbf{Y}^t$  and  $\mathbf{Q}_0$ ); i.e.

$$\mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) = \begin{bmatrix} q_t^1 \pi_{11} \psi_1(\mathbf{y}^*) + q_t^2 \pi_{21} \psi_2(\mathbf{y}^*) + \cdots + q_t^n \pi_{n1} \psi_n(\mathbf{y}^*) \\ q_t^1 \pi_{12} \psi_1(\mathbf{y}^*) + q_t^2 \pi_{22} \psi_2(\mathbf{y}^*) + \cdots + q_t^n \pi_{n2} \psi_n(\mathbf{y}^*) \\ \vdots \\ q_t^1 \pi_{1n} \psi_1(\mathbf{y}^*) + q_t^2 \pi_{2n} \psi_2(\mathbf{y}^*) + \cdots + q_t^n \pi_{nn} \psi_n(\mathbf{y}^*) \end{bmatrix}_{n \times 1}.$$

We wish to write this more succinctly.

First, observe that, if  $\mathbf{X}_t$  is the  $i$ th coordinate vector, then  $\mathbf{X}_t \mathbf{X}_t'$  is an  $n \times n$  matrix with 1 on the  $i$ th diagonal and zero everywhere else. Since  $\mathbf{Q}_t$  is the vector of (conditional) probabilities that the current state is the 1st, 2nd, ...,  $n$ th state,

$$\begin{aligned} \mathbb{E}[\mathbf{X}_t \mathbf{X}_t' | \mathbf{Q}_t] &= q_t^1 \text{diag}[\mathbf{e}_1] + q_t^2 \text{diag}[\mathbf{e}_2] + \cdots + q_t^n \text{diag}[\mathbf{e}_n] \\ &= \text{diag}[q_t^1 \mathbf{e}_1] + \text{diag}[q_t^2 \mathbf{e}_2] + \cdots + \text{diag}[q_t^n \mathbf{e}_n] \\ &= \text{diag}[\mathbf{Q}_t], \end{aligned}$$

where  $\mathbf{e}_i$  is the  $i$ th coordinate vector and where  $\text{diag}[\mathbf{Q}_t]$  is a diagonal matrix with the entries of  $\mathbf{Q}_t$  on the diagonal. The joint distribution that we want can therefore be expressed as

$$\begin{aligned} \mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t, \mathbf{Y}^t, \mathbf{Q}_0) &= \mathbb{E}[\Pi' \mathbf{X}_t \mathbf{X}_t' \text{vec}[\psi_i(\mathbf{y}^*)] | \mathbf{Q}_t] \\ &= \Pi' \mathbb{E}[\mathbf{X}_t \mathbf{X}_t' | \mathbf{Q}_t] \text{vec}[\psi_i(\mathbf{y}^*)] \\ &= \Pi' \text{diag}[\mathbf{Q}_t] \text{vec}[\psi_i(\mathbf{y}^*)], \end{aligned}$$

Thus,  $\mathbf{Q}_t$  encodes the information in history of signals that is relevant for this calculation. Observe that, conditioned on  $\mathbf{Q}_t$ , the distributions for  $\mathbf{X}_{t+1}$  and  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  are not independent.<sup>4</sup>

**Step 3: Find the implied distribution for  $\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t$ .** By law of total probabilities:

$$\mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) = \sum_{i=1}^n \mathbb{P}(\mathbf{X}_{t+1} = \mathbf{e}_i, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t).$$

That is, to obtain the distribution of  $\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t$ , we need to sum across the rows (or elements) of (??), which is equivalent to summing over the “hidden” states:

$$\begin{aligned} \mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) &= q_t^1 (\pi_{11} + \pi_{12} + \cdots + \pi_{1n}) \psi_1(\mathbf{y}^*) + q_t^2 (\pi_{21} + \pi_{22} + \cdots + \pi_{2n}) \psi_2(\mathbf{y}^*) \\ &\quad + \cdots + q_t^n (\pi_{n1} + \pi_{n2} + \cdots + \pi_{nn}) \psi_n(\mathbf{y}^*) \\ &= q_t^1 \psi_1(\mathbf{y}^*) + q_t^2 \psi_2(\mathbf{y}^*) + \cdots + q_t^n \psi_n(\mathbf{y}^*) \\ &= \underbrace{\mathbf{1}_n' \Pi' \text{diag}[\mathbf{Q}_t]}_{=\mathbf{Q}_t'} \text{vec}[\psi_i(\mathbf{y}^*)], \end{aligned}$$

---

<sup>4</sup>Observe that

$$\begin{aligned} \mathbb{P}(\mathbf{X}_{t+1} | \mathbf{Q}_t) &= \Pi' \mathbf{Q}_t, \\ \mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) &= \mathbf{Q}_t' \text{vec}[\psi_i(\mathbf{y}^*)], \\ \mathbb{P}(\mathbf{X}_{t+1} | \mathbf{Q}_t) \mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) &= \Pi' \mathbf{Q}_t \mathbf{Q}_t' \text{vec}[\psi_i(\mathbf{y}^*)] \\ &\neq \Pi' \text{diag}[\mathbf{Q}_t] \text{vec}[\psi_i(\mathbf{y}^*)]. \end{aligned}$$

where we used the fact that

$$\mathbf{1}'_n \Pi' = \mathbf{1}'_n \Rightarrow \mathbf{1}'_n \Pi' \text{diag}[\mathbf{Q}_t] = [q_t^1, q_t^2, \dots, q_t^n] = \mathbf{Q}'_t.$$

Observe that  $\mathbf{Q}_t$  is a vector of weight used in forming a mixture distribution. Suppose, for instance,  $\psi_i$  is a normal distribution with mean  $\boldsymbol{\mu}_i$  and variance-covariance matrix  $\boldsymbol{\Sigma}_i$ , then the distribution of  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  conditioned on  $\mathbf{Q}_t$  is a mixture of normals with mixing probabilities given by respective entries of  $\mathbf{Q}_t$

**Step 4: Compute  $\mathbf{Q}_{t+1}$ .** Substituting the expressions we obtained into the equation for  $\mathbf{Q}_{t+1}$  yields

$$\begin{aligned} \mathbf{Q}_{t+1} &= \frac{\mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)}{\mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)} \\ &= \frac{\Pi' \text{diag}(\mathbf{Q}_t) \text{vec}[\psi_i(\mathbf{Y}_{t+1} - \mathbf{Y}_t)]}{\mathbf{Q}'_t \text{vec}[\psi_i(\mathbf{Y}_{t+1} - \mathbf{Y}_t)]}. \end{aligned}$$

Observe that steps (3) and (4) define a Markov process for  $\mathbf{Q}_{t+1}$ . Step (3) gives  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  drawn from a (history dependent) mixture of densities  $\psi_i$ , and step (4), we construct  $\mathbf{Q}_{t+1}$  as a function of  $\mathbf{Q}_t$  and  $\mathbf{Y}_{t+1} - \mathbf{Y}_t$  from step (3).

**Algorithm 2.** (Learning about discrete states).  $\{\mathbf{X}_t\}$  is a hidden  $n$ -state Markov process with (a known) transition matrix  $\Pi$ . There is a vector of signals denoted by  $\{\mathbf{Y}_{t+1} - \mathbf{Y}_t\}$  with density  $\psi_i(\mathbf{y}^*)$  if the state is  $i$  (i.e.  $\mathbf{X}_t$  is the  $i$ th coordinate vector).  $\mathbf{Q}_t$  denotes the probability that the state is  $i$  given the signal history, and  $\mathbf{Q}_0$  is the vector of initial probabilities.  $\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{X}_t$  and  $\mathbf{X}_{t+1} | \mathbf{X}_t$  are independent.

▷ Step 1: Find the joint distribution for  $(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t) | \mathbf{X}_t$

$$\begin{aligned} \mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{X}_t) &= \mathbb{P}(\mathbf{X}_{t+1} | \mathbf{X}_t) \mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{X}_t) \\ &= \Pi' \mathbf{X}_t \mathbf{X}'_t \text{vec}[\psi_i(\mathbf{y}^*)]. \end{aligned}$$

▷ Step 2: Find the joint distribution of  $(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t) | \mathbf{Q}_t$  by taking the expectation of the above term.

$$\begin{aligned} \mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) &= \mathbb{E}[\Pi' \mathbf{X}_t \mathbf{X}'_t \text{vec}[\psi_i(\mathbf{y}^*)] | \mathbf{Q}_t] \\ &= \Pi' \text{diag}[\mathbf{Q}_t] \text{vec}[\psi_i(\mathbf{y}^*)]. \end{aligned}$$

▷ Step 3: Find the implied distribution for  $\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t$ .

$$\begin{aligned} \mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) &= \sum_{i=1}^n \mathbb{P}(\mathbf{X}_{t+1} = \mathbf{x}_i, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t) \\ &= \underbrace{\mathbf{1}'_n \Pi' \text{diag}[\mathbf{Q}_t]}_{=\mathbf{Q}'_t} \text{vec}[\psi_i(\mathbf{y}^*)], \end{aligned}$$

▷ Step 4: Compute  $\mathbf{Q}_{t+1}$ .

$$\begin{aligned} \mathbf{Q}_{t+1} &= \frac{\mathbb{P}(\mathbf{X}_{t+1}, \mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)}{\mathbb{P}(\mathbf{Y}_{t+1} - \mathbf{Y}_t | \mathbf{Q}_t)} \\ &= \frac{\Pi' \text{diag}(\mathbf{Q}_t) \text{vec}[\psi_i(\mathbf{Y}_{t+1} - \mathbf{Y}_t)]}{\mathbf{Q}'_t \text{vec}[\psi_i(\mathbf{Y}_{t+1} - \mathbf{Y}_t)]}. \end{aligned}$$

## 7 Generalized Method of Moments

Generalised Method of Moments (GMM) refers to a class of estimators which are constructed from exploiting the sample moment counterparts of population moment conditions (sometimes known as orthogonality conditions) of the data generating model. GMM estimators have become widely used, for the following reasons.

- ▷ GMM estimators have large sample properties that are easy to characterise in ways that facilitate comparison. A family of such estimators can be studied a priori in ways that make asymptotic efficiency comparisons easy. The method also provides a natural way to construct tests which take account of both sampling and estimation error.
- ▷ In practice, researchers find it useful that GMM estimators can be constructed without specifying the full data generating process (which would be required to write down the maximum likelihood estimator.) This characteristic has been exploited in analysing partially specified economic models, in studying potentially misspecified dynamic models designed to match target moments, and in constructing stochastic discount factor models that link asset pricing to sources of macroeconomic risk.

### 7.1 Preliminaries

Let  $\mathbf{Z}$  be an index set of  $\mathcal{Z}$ , where  $\mathcal{Z}$  is a linear space of  $m$ -dimensional random vectors. Let  $\Phi$  be a known  $m$ -dimensional function and suppose that an economic model implies

$$\mathbb{E}[\mathbf{Z} \cdot \Phi(\mathbf{Y}, \boldsymbol{\beta})] = 0, \forall \mathbf{Z} \in \mathcal{Z},$$

where  $\boldsymbol{\beta} \in \mathbb{R}^k$  is an unknown parameter and  $\mathbf{Y}$  is the data. The function  $\Phi$  can be non-linear.

Given that  $\boldsymbol{\beta}$  is a  $k \times 1$  vector, we require  $\mathbf{Z}^j$  entries in  $\mathcal{Z}$  for  $j = 1, 2, \dots, k$  which gives us  $k$  equations in  $k$  unknowns:

$$\mathbb{E}[\mathbf{Z}^j \cdot \Phi(\mathbf{Y}, \mathbf{b})] = 0, j = 1, 2, \dots, k$$

to use in identifying  $\mathbf{b} = \boldsymbol{\beta}$ . We can think of the  $\mathbf{Z}^j$ 's as "instruments". The question then becomes how should we choose  $\mathbf{Z}^j$ 's and what are the properties of the estimator for  $\boldsymbol{\beta}$  obtained from this procedure?

#### 7.1.1 Finite unconditional moment restrictions

The set  $\mathcal{Z}$  can potentially be very large. Here, we collapse  $\mathcal{Z}$  so that it consists of linear combinations of  $r \geq k$  basis  $\mathbf{Z}$ 's that satisfy  $\mathbb{E}[\mathbf{Z} \cdot \Phi(\mathbf{Y}, \boldsymbol{\beta})] = 0, \forall \mathbf{Z} \in \mathcal{Z}$ . Let

$$F(\mathbf{X}, \mathbf{b})_{r \times 1} = \begin{pmatrix} \mathbf{Z}^{1'} \\ \mathbf{Z}^{2'} \\ \vdots \\ \mathbf{Z}^{r'} \end{pmatrix}_{r \times m} \Phi(\mathbf{X}, \mathbf{b})_{m \times 1}.$$

The unconditional moment restriction is

$$\mathbb{E}[F(\mathbf{X}, \boldsymbol{\beta})] = \mathbf{0}_{r \times 1},$$

where  $\mathbf{X}$  contains  $\mathbf{Y}$  and the  $r$  basis  $\mathbf{Z}$ 's. Let  $\mathcal{Z} = \mathbb{R}^r$ . In this case, an estimator may be associated with an  $r \times k$  selection matrix  $\mathbf{A}$  or real numbers where each of the  $k$  columns are vectors in  $\mathcal{Z}$ . Thus, the equation used to estimate are

$$\mathbf{A}' \mathbb{E}[F(\mathbf{X}, \mathbf{b})] = \mathbf{0}_{k \times 1},$$

which are satisfied for  $\mathbf{b} = \boldsymbol{\beta}$ . A choice of  $\mathbf{A}$  implements/indexes a particular GMM estimator.

**Example 7.1.** (Moment matching) Suppose that

$$\Phi(\mathbf{Y}, \mathbf{b}) = \Psi(\mathbf{Y}) - \Gamma(\mathbf{b}),$$

where

$$\mathbb{E}[\Psi(\mathbf{Y})] = \Gamma(\beta).$$

In this formulation,  $\Psi(\mathbf{Y})$  defines the moments to be matched and  $\Gamma(\mathbf{b})$  give the predicted moments from a model as a function of a potential parameter vector  $\mathbf{b}$ . We presume that  $r \geq k$  and  $\mathcal{Z} = \mathbb{R}^r$ . In other words, LHS can be obtained from data and the RHS can be obtained from the model.

**Example 7.2.** (IV) Suppose  $\alpha \cdot Y_t = u_t$  (disturbances) and suppose I know that  $\mathbb{E}[Z_t u_t] = 0$  where  $Z_t$  is a  $r \times 1$  matrix. This is equivalent to saying:

$$\alpha' \mathbb{E}[Y_t Z_t'] = \mathbf{0}$$

In the population,  $\alpha$  is only identified up to a scalar. So you need  $\beta$  to impose conditions on  $\alpha$ . What you're identifying here is a null-space and you use  $\beta$  to get identification.

### 7.1.2 Conditional moment restrictions

Consider the following conditional moment restriction:

$$\mathbb{E}[\Phi(\mathbf{Y}, \beta) | \mathcal{K}] = 0,$$

where  $\mathcal{K}$  is some conditioning information set.

Suppose that  $\Phi(\mathbf{Y}, \mathbf{b})$  for  $\mathbf{b} \in \mathbb{R}^k$  has a finite second moment and  $\mathcal{Z}$  contains at least  $m$ -dimensional random vectors that are bounded and  $\mathcal{K}$  measurable.

By Law of Iterated Expectations, we can see that conditional moment restriction implies unconditional moment restriction:

$$\mathbb{E}[\mathbf{Z} \cdot \Phi(\mathbf{Y}, \beta)] = \mathbb{E}[\mathbf{Z} \cdot \mathbb{E}[\Phi(\mathbf{Y}, \beta) | \mathcal{K}]] = 0$$

since  $\mathbf{Z}$  is in the  $\mathcal{K}$  information set. In general, the converse is not true; i.e. unconditional moment restriction does not imply conditional moment restriction.

However, if the unconditional moment restriction holds for all  $\mathbf{Z} \in \mathcal{Z}$ , then we have the orthogonality condition for the conditional expectation:

$$\mathbb{E}[\mathbf{Z} \cdot (\Phi(\mathbf{Y}, \beta) - \mathbf{0})] = 0, \forall \mathbf{Z} \in \mathcal{Z}$$

so that we can get back to

$$\mathbb{E}[\Phi(\mathbf{Y}, \beta) | \mathcal{K}] = 0.$$

So, in this case, the family of unconditional moment restrictions imply the conditional moment restriction.

**Example 7.3.** (Stochastic discount factor) Consider the following moment restriction

$$\mathbb{E}\left[\left(\frac{S_{t+\ell}}{S_t}\right) \mathbf{R}_{t+\ell} | \mathfrak{F}_t\right] = \mathbf{1}_n,$$

where  $S_{t+\ell}/S_t$  is a stochastic discount factor over horizon  $\ell$ ,  $\mathbf{R}_{t+\ell}$  is an  $n$ -dimensional vector of returns, and  $\mathfrak{F}_t$  is the information available at time  $t$ . Here,

$$\begin{aligned} \Phi(\mathbf{Y}, \mathbf{b}) &= \left(\frac{S_{t+\ell}}{S_t}\right) \mathbf{R}_{t+\ell} - \mathbf{1}_n, \\ \mathcal{K} &= \mathfrak{F}_t. \end{aligned}$$

## 7.2 Traditional Way of Presenting GMM (via Minimization Problem)

Given some parameter space  $\mathbf{P}$ , we can obtain the GMM estimator  $\mathbf{b}_N$  by solving

$$\min_{\mathbf{b} \in \mathbf{P}} \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \right)' (\mathbf{V}_N(\mathbf{b}))^{-1} \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \right).$$

Notice that we are simultaneously solving for the estimator  $\mathbf{b}$  as well as the variance  $\mathbf{V}_N(\mathbf{b})$ .

The first-order condition is

$$\mathbf{0} = 2 \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N \frac{\partial}{\partial \mathbf{b}} F(\mathbf{X}_t, \mathbf{b}) \right)' (\mathbf{V}_N(\mathbf{b}))^{-1} \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \right) + \text{other terms}.$$

(The other terms, such as the derivative of  $(\mathbf{V}_N(\mathbf{b}))^{-1}$ , are going to be pre- and post- multiplied by  $N^{-1/2} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b})$ , which converges in probability to zero so we need not worry about them). Notice that

$$\frac{1}{N} \sum_{t=1}^N \frac{\partial}{\partial \mathbf{b}} F(\mathbf{X}_t, \mathbf{b}) \xrightarrow{P} \mathbf{D} = \mathbb{E} \left[ \frac{\partial F(\mathbf{X}, \boldsymbol{\beta})}{\partial \mathbf{b}} \right]$$

So asymptotically,

$$\left( \frac{1}{N} \sum_{t=1}^N \frac{\partial}{\partial \mathbf{b}} F(\mathbf{X}_t, \mathbf{b}) \right)' (\mathbf{V}_N(\mathbf{b}))^{-1} \xrightarrow{P} \mathbf{D}' \mathbf{V}^{-1} = \mathbf{A}'.$$

The first-order condition in the limit is

$$\mathbf{0} \approx \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}).$$

Thus, the estimator is asymptotically efficient.

## 7.3 GMM Estimation with Constant Selection

Consider  $F(\mathbf{x}, \mathbf{b}_{k \times 1})_{r \times 1}$  and consider  $\mathbb{E}[F(\mathbf{x}, \mathbf{b}) | \ell] = 0$  (where we will drop conditioning). Then consider

$$\mathbf{A}_{r \times 1}^i \cdot \mathbb{E}[F(\mathbf{X}_t, \mathbf{b})] = \mathbf{0}, \forall t = 1, \dots, k$$

so we have  $k$  equations yielding a unique  $\mathbf{b}_N$ . Stacking  $\mathbf{A} = [\mathbf{A}^1, \dots, \mathbf{A}^k]$ , we have

$$\mathbf{A}' \mathbb{E}[F(\mathbf{X}_t, \mathbf{b})] = \mathbf{0} \Leftrightarrow \mathbf{b} = \boldsymbol{\beta}$$

Then we say  $\mathbf{b}_N$  is a GMM estimator where

$$\mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N) = \mathbf{0}$$

where  $\mathbf{A}$  is the selection matrix. We are interested in studying a family of estimators indexed by  $\mathbf{A}$ .

### 7.3.1 Example: Euler Equations and SDF

Consider a consumption process  $\{C_t\} : (C_0, C_1, \dots)$ . We can think of this being on an infinite-dimensional indifference curve. Now stay on same indifference curve and consider a different point for some scalar  $r$ :

$$\left( \underline{C_0 - P_0(r)}, C_1, C_2, \dots, \underline{C_\tau + r\xi_\tau}, C_{\tau+1}, \dots \right)$$

where  $\xi_\tau$  is a random variable dependent on time  $\tau$ . To change time- $\tau$  consumption by that amount, we need to reduce today's consumption by  $P_0(r)$ .

Now define  $\pi_0^\tau(S_\tau)$ , the slope of the indifference curve, to be

$$\pi_0^\tau(\xi_\tau) = \frac{d}{dr} P_0(r) |_{r=0}$$

and this also gives us the shadow price of asset (security) that pays off  $\xi_\tau$  at date  $\tau$ . Using a SDF notation:

$$\pi_0^\tau(\xi_\tau) = \mathbb{E} \left[ \left( \frac{S_\tau}{S_0} \right) \xi_\tau | \mathfrak{F}_0 \right]$$

where  $S_{t+\ell}/S_t$  is a stochastic discount factor over horizon  $\tau$ .

- ▷ Stochastic: Since  $\xi_\tau$  is random, we can't discount it at some fixed rate. The riskiness of  $\xi_\tau$  has to demand some kind of compensation. The way you discount the cash flow by something that is random – you discount different states of the world differently.

So how should we get the SDF? We know that prices equal the marginal rates of substitution, so we use this insight:

$$\frac{S_\tau}{S_0} = \frac{MU_\tau^0}{MU_0^0}$$

where the RHS is the intertemporal marginal rate of substitution. To proceed, we need to define the preferences (from the perspective from time zero) as

$$\mathbb{E} \left[ \sum_{t=0}^{\infty} \exp(-\delta\tau) U(C_t) | \mathfrak{F}_0 \right]$$

Then  $MU_\tau^0 = \exp(-\delta\tau) U'(C_\tau)$  and  $MU_0^0 = U'(C_0)$  which gives us:

$$\frac{S_\tau}{S_0} = \exp(-\delta\tau) \frac{U'(C_\tau)}{U'(C_0)}$$

Thus:

$$\pi_0^\tau(\xi_\tau) = \mathbb{E} \left[ \exp(-\delta\tau) \frac{U'(C_\tau)}{U'(C_0)} \xi_\tau | \mathfrak{F}_0 \right]$$

This is called the Euler Equation since we are sitting in the optimal consumption process and throwing in this perturbation to see how consumption changes. This is an example of the conditional moment restriction.

### 7.3.2 Adding Distributional Assumptions

Now for pedagogical reasons, we add some distributional assumptions. Note that these conditions are not necessary for the estimation. Assume  $\xi_\tau$  is the gross return i.e, an object that has price  $\pi_0^\tau(\xi_\tau) = 1$  and  $(\xi_\tau, S_\tau/S_0)$  are jointly lognormal. Furthermore:

$$\begin{aligned} \log S_\tau - \log S_0 &= \mu_0^S + \sigma^S \cdot W_\tau, & W_\tau &\sim N(\mathbf{0}, I) \\ \log \xi_\tau &= \mu_0^\xi + \sigma^\xi \cdot W_\tau \end{aligned}$$

Note that  $\mu_0^\xi$  may also depend on  $\tau$ .

▷ Adding the two equations, we have:

$$\begin{aligned}\log S_\tau - \log S_0 + \log \xi_\tau &= \left(\mu_0^s + \mu_0^\xi\right) + \left(\sigma^S + \sigma^\xi\right) W_\tau \\ \Rightarrow \log \frac{S_\tau}{S_0} \xi_\tau &= \left(\mu_0^s + \mu_0^\xi\right) + \left(\sigma^S + \sigma^\xi\right) W_\tau\end{aligned}$$

▷ Therefore:

$$\log \mathbb{E} \left[ \frac{S_\tau}{S_0} \xi_\tau | \mathfrak{F}_0 \right] = \mu_0^s + \mu_0^\xi + \frac{1}{2} \left| \sigma^S + \sigma^\xi \right|^2$$

which must be equal to zero since  $\pi_0^\tau(\xi_\tau) = 1$ .

▷ Rearranging:

$$\begin{aligned}0 &= \mu_0^s + \mu_0^\xi + \frac{1}{2} (\sigma^S)^2 + \frac{1}{2} (\sigma^\xi)^2 + \left| \sigma^S \cdot \sigma^\xi \right| \\ \Leftrightarrow -\sigma^S \cdot \sigma^\xi &= \underbrace{\left[ \mu_0^\xi + \frac{1}{2} |\sigma^\xi|^2 \right]}_{[1]} + \underbrace{\left[ \mu_0^s + \frac{1}{2} |\sigma^S|^2 \right]}_{[2]} \quad [A]\end{aligned}$$

\* [1] is equal to the log expected return, and [2] is minus of log riskless return. To get the riskless return, set  $\sigma^\xi = 0$ , then the (riskfree) return is  $\mu_0^\xi$ . We can explicitly obtain this expression from [A]:

$$0 = \mu_0^\xi - \left( \mu_0^s + \frac{1}{2} |\sigma^S|^2 \right) \Rightarrow \mu_0^\xi = - \left( \mu_0^s + \frac{1}{2} |\sigma^S|^2 \right)$$

\* Therefore,  $-\sigma^S \cdot \sigma^\xi$  is the risk compensation. We can think of  $-\sigma^S$  as the risk price vector. It is linear in terms of how exposed you are to the risk.

## 7.4 GMM Limiting Approximation

Start with

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) \xrightarrow{d} N(0, \mathbf{V}), \quad \mathbf{V} := \mathbb{E} [\mathbf{H}_{t+1} \mathbf{H}_{t+1}']$$

How can we get here? Take consistency of  $\mathbf{b}_N$  as given:

▷ Taylor approximation of  $F$  around  $\mathbf{b}_N = \beta$  gives

$$\mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N) \approx \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) + \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N \frac{\partial F(\mathbf{X}, \beta)}{\partial \mathbf{b}'} (\mathbf{b}_N - \beta).$$

You can also get the same result using the mean-value theorem, but this introduces a new point  $\bar{\mathbf{b}}_N$  which will also converge to  $\beta$ .

▷ Mean-value theorem with law of large numbers implies

$$\frac{1}{N} \sum_{t=1}^N \frac{\partial F(\mathbf{X}, \beta)}{\partial \mathbf{b}'} \xrightarrow{p} \mathbb{E} \left[ \frac{\partial F(\mathbf{X}, \beta)}{\partial \mathbf{b}'} \right] =: \mathbf{D},$$



where  $\mathbf{D}$  is an  $r \times k$  vector of expectation of partial derivatives of  $F$  with respect to  $\mathbf{b}$ ; i.e.

$$\mathbf{D} := \mathbb{E} \begin{bmatrix} \frac{\partial F_1(\mathbf{X}, \boldsymbol{\beta})}{\partial b_1} & \frac{\partial F_1(\mathbf{X}, \boldsymbol{\beta})}{\partial b_2} & \dots & \frac{\partial F_1(\mathbf{X}, \boldsymbol{\beta})}{\partial b_k} \\ \frac{\partial F_2(\mathbf{X}, \boldsymbol{\beta})}{\partial b_1} & \frac{\partial F_2(\mathbf{X}, \boldsymbol{\beta})}{\partial b_2} & \dots & \frac{\partial F_2(\mathbf{X}, \boldsymbol{\beta})}{\partial b_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_r(\mathbf{X}, \boldsymbol{\beta})}{\partial b_1} & \frac{\partial F_r(\mathbf{X}, \boldsymbol{\beta})}{\partial b_2} & \dots & \frac{\partial F_r(\mathbf{X}, \boldsymbol{\beta})}{\partial b_k} \end{bmatrix}_{r \times k}.$$

So,

$$\underbrace{\mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N)}_{k \text{ equations stacked}} \approx \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}) + \mathbf{A}' \sqrt{N} \frac{1}{N} \sum_{t=1}^N \frac{\partial F(\mathbf{X}, \boldsymbol{\beta})}{\partial \mathbf{b}'} (\mathbf{b}_N - \boldsymbol{\beta})$$

$$\approx \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}) + \mathbf{A}' \mathbf{D} \sqrt{N} (\mathbf{b}_N - \boldsymbol{\beta}).$$

▷ Since the left-hand side equals zero given that we are considering estimators of the kind

$$\mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N) = \mathbf{0}$$

rearranging gives

$$\mathbf{A}' \mathbf{D} \sqrt{N} (\mathbf{b}_N - \boldsymbol{\beta}) \approx \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}).$$

Assuming that  $\mathbf{A}' \mathbf{D}$  is nonsingular, then

$$\sqrt{N} (\mathbf{b}_n - \boldsymbol{\beta}) \approx - (\mathbf{A}' \mathbf{D})^{-1} \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}).$$

The key here is that the covariance matrix depends on  $\mathbf{A}$  and  $\mathbf{D}$  so the natural question is bounding this covariance.

## 7.5 GMM Efficiency Bound

Every GMM estimator satisfies:

$$\sqrt{N} (\mathbf{b}_n - \boldsymbol{\beta}) \approx - (\mathbf{A}' \mathbf{D})^{-1} \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}),$$

We also know from the Central Limit Theorem:

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta}) \xrightarrow{d} N(0, \mathbf{V}), \quad \mathbf{V} = \mathbb{E} [\mathbf{H}_{t+1} \mathbf{H}_{t+1}'].$$

Combining these two together, we have:

$$\sqrt{N} (\mathbf{b}_n - \boldsymbol{\beta}) \xrightarrow{d} N(0, \text{Cov}[\mathbf{A}]),$$

where  $\text{Cov}[\mathbf{A}]$  denotes the asymptotic variance-covariance of the estimator (indexed by the selection matrix  $\mathbf{A}$ ):

$$\text{Cov}[\mathbf{A}] := (\mathbf{A}' \mathbf{D})^{-1} \mathbf{A}' \mathbf{V} \mathbf{A} (\mathbf{D}' \mathbf{A})^{-1}.$$

The goal is to find the greatest lower bound for  $\text{Cov}[\mathbf{A}]$  provided that there exist an  $\mathbf{A}$  such that  $\mathbf{A}' \mathbf{D}$  is non-singular.

1. Establish that post-multiplying the weighting matrix with a nonsingular matrix does not alter the asymptotic variance-covariance matrix of the GMM estimator.

▷ This is intuitive because premultiplying the estimation equation

$$\mathbf{A}' \frac{1}{N} \sum_{t=1}^N F_t(\mathbf{X}_t, \mathbf{b}_N) = \mathbf{0}.$$

will not change the optimal  $\mathbf{A}$ .

2. Find  $\tilde{\mathbf{A}}$  such that  $\mathbf{A}'\mathbf{V}\tilde{\mathbf{A}} = \mathbf{A}'\mathbf{D}$  for all  $\mathbf{A}$ . It turns out to be  $\tilde{\mathbf{A}} = \mathbf{V}^{-1}\mathbf{D}$ .

▷ This is a step to help us figure out  $\mathbf{A}^*$  in the next step.

3. Above result allows us to normalize  $\mathbf{A}$  to  $\mathbf{A}^*$  such that  $\langle \mathbf{A}^* | \tilde{\mathbf{A}} \rangle = \mathbf{A}^*\mathbf{D} = \mathbf{I}$  can be done without loss of efficiency, and without affecting the GMM estimators.

4. Denote  $\mathbf{A}^* = \tilde{\mathbf{A}} (\mathbf{D}'\tilde{\mathbf{A}})^{-1} = \mathbf{V}^{-1}\mathbf{D} (\mathbf{D}'\mathbf{V}^{-1}\mathbf{D})^{-1}$  and show that  $Cov[\mathbf{A}^*]$  is the GMM efficiency bound, i.e.

$$Cov[\mathbf{A}] \geq Cov[\mathbf{A}^*] = Cov[\tilde{\mathbf{A}}]$$

▷ First, observe that  $\mathbf{A}'\mathbf{V}\mathbf{A}^* = \mathbf{A}'\mathbf{D} (\mathbf{D}'\mathbf{V}^{-1}\mathbf{D})^{-1} = (\mathbf{D}\mathbf{V}^{-1}\mathbf{D})^{-1}$  provided  $\mathbf{A}'\mathbf{D} = \mathbf{I}$ .

▷ Second, note that

$$\begin{aligned} Cov[\mathbf{A}] &= (\mathbf{A}'\mathbf{V}\tilde{\mathbf{A}})^{-1} \mathbf{A}'\mathbf{V}\mathbf{A} (\tilde{\mathbf{A}}'\mathbf{V}\mathbf{A})^{-1} \\ &= \mathbf{I}^{-1} \mathbf{A}'\mathbf{V}\mathbf{A} \mathbf{I} = \mathbf{A}'\mathbf{V}\mathbf{A} \end{aligned}$$

▷ Then

$$\begin{aligned} (\mathbf{A} - \mathbf{A}^*)' \mathbf{V} (\mathbf{A} - \mathbf{A}^*) &= \mathbf{A}'\mathbf{V}\mathbf{A} - \mathbf{A}^*\mathbf{V}\mathbf{A} - \mathbf{A}'\mathbf{V}\mathbf{A}^* + \mathbf{A}^{*'}\mathbf{V}\mathbf{A}^* \\ &= \mathbf{A}'\mathbf{V}\mathbf{A} - \mathbf{A}^*\mathbf{V}\mathbf{A}^* \\ &= Cov[\mathbf{A}] - Cov[\mathbf{A}^*] \end{aligned}$$

where the penultimate equality is from the fact that  $\mathbf{A}^{*'}\mathbf{D} = \mathbf{I}$  and thus we can safely plug in  $\mathbf{A}^*$  into the equation in the first bullet.

▷ The LHS is positive semi-definite by construction so the RHS must also be positive semi-definite.

Just remember that  $\mathbf{A}^* = \mathbf{V}^{-1}\mathbf{D}$  is the efficient selection matrix where  $\mathbf{D}_{r \times k} = \mathbb{E}[\partial F(\mathbf{X}_t, \beta) / \partial \mathbf{b}']$ .

**Example 7.4.** (Imposing Restrictions) Consider  $\alpha' \mathbb{E}[Y_t | \mathcal{F}_{t-\ell}] = 0, \ell \geq 1$ . Suppose we have  $z_{t-\ell} \in \mathcal{F}_{t-\ell}$ . Then we have a matrix

$$\mathbb{E}[z_{t-\ell} Y_t'] \alpha = 0$$

where  $\beta$  restricts  $\alpha$ . Regress  $y_t^1$  onto  $z_{t-\ell}$  and regress  $y_t^2$  onto  $z_{t-\ell}$ .

▷ If both of them have zero predictability, we are in trouble because  $\mathbb{E}[z_{t-\ell} Y_t'] \alpha$  will be all zeros.

▷ I just need one of these two regression to tell me that it's predictable.

Think of the structural model as  $\alpha' y_t = u_t$  and we need  $\mathbb{E}[z_{t-\ell} u_t] = 0$  for things to be working. If it has rank two, we're fucked. If it has rank one, we're good but to make it rank one, we can have:

$$\alpha = \begin{bmatrix} 1 \\ \beta \end{bmatrix} \quad vs. \quad \alpha = \begin{bmatrix} \beta \\ 1 \end{bmatrix}$$

which yields different

$$\mathbb{E} \left[ \frac{\partial F(\mathbf{X}, \beta)}{\partial \mathbf{b}'} \right] =: \mathbf{D} = \mathbb{E} [y_t^2 z_{t-2}] \quad vs. \quad \mathbb{E} [y_t^1 z_{t-2}]$$

**Example 7.5.** (*Serial correlation*) Suppose we have quarterly data and the moment condition as

$$\mathbb{E} [F(\mathbf{X}_t, \beta) | \mathfrak{F}_{t-\ell}] = \mathbf{0}.$$

Furthermore, denote  $G_t = F(\mathbf{X}_t, \beta)$ . Let  $\ell = 2$ . Then,

$$\log Y_{t+2} - \log Y_t = \text{constant} + \text{parameters} + u_{t+2},$$

where  $u_{t+2} = F(\mathbf{X}_t, \beta)$ . We assume

$$\mathbb{E}[u_{t+2} | \mathfrak{F}_t] = \mathbb{E}[F(\mathbf{X}_{t+2}, \beta) | \mathfrak{F}_t] = 0.$$

By law of iterated expectations,

$$\begin{aligned} \mathbb{E} [F(\mathbf{X}_{t+2}, \beta) F(\mathbf{X}_t, \beta)'] &= \mathbb{E} \left[ \underbrace{\mathbb{E} [F(\mathbf{X}_{t+2}, \beta) | \mathfrak{F}_t]}_{=0} F(\mathbf{X}_t, \beta)' \right] = 0, \\ \mathbb{E} [F(\mathbf{X}_{t+3}, \beta) F(\mathbf{X}_t, \beta)'] &= \mathbb{E} \left[ \underbrace{\mathbb{E} [F(\mathbf{X}_{t+3}, \beta) | \mathfrak{F}_{t+1}]}_{=0} F(\mathbf{X}_t, \beta)' \right] = 0, \\ \mathbb{E} [F(\mathbf{X}_{t-2}, \beta) F(\mathbf{X}_t, \beta)'] &= \mathbb{E} \left[ F(\mathbf{X}_{t-2}, \beta) \underbrace{\mathbb{E} [F(\mathbf{X}_t, \beta)' | \mathfrak{F}_{t-2}]}_{=0} \right] = 0, \\ \mathbb{E} [F(\mathbf{X}_{t-3}, \beta) F(\mathbf{X}_t, \beta)'] &= \mathbb{E} \left[ F(\mathbf{X}_{t-3}, \beta) \underbrace{\mathbb{E} [F(\mathbf{X}_t, \beta)' | \mathfrak{F}_{t-2}]}_{=0} \right] = 0. \end{aligned}$$

So, any cross products with time subscript difference greater or equal to two is zero. But, note that we have not imposed

restrictions on  $\mathbb{E} [F(\mathbf{X}_{t+1}, \beta) F(\mathbf{X}_t, \beta)']$ ,  $\mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_t, \beta)']$ , or  $\mathbb{E} [F(\mathbf{X}_{t-1}, \beta) F(\mathbf{X}_t, \beta)']$ . This means that

$$\begin{aligned} \mathbf{V} &= \lim_{N \rightarrow \infty} \mathbb{E} \left[ \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) \right) \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) \right)' \right] \\ &= \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{1}{N} \sum_{t=1}^N F(\mathbf{X}_t, \beta) F(\mathbf{X}_t, \beta)' \right] + \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{1}{N} \sum_{t=1}^{N-1} F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t+1}, \beta)' \right] \\ &\quad + \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{1}{N} \sum_{t=2}^N F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t-1}, \beta)' \right] \\ &= \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{1}{N} \sum_{t=1}^N F(\mathbf{X}_t, \beta) F(\mathbf{X}_t, \beta)' \right] + \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{N-1}{N} \frac{1}{N-1} \sum_{t=1}^{N-1} F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t+1}, \beta)' \right] \\ &\quad + \lim_{N \rightarrow \infty} \mathbb{E} \left[ \frac{N-1}{N} \frac{1}{N-1} \sum_{t=2}^N F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t-1}, \beta)' \right]. \end{aligned}$$

Equivalently, we can visualise above as:

$$\begin{aligned} &\mathbb{E} [F(\mathbf{X}_t, \beta) (F(\mathbf{X}_j, \beta))'] \\ &= \begin{bmatrix} \mathbb{E} [F_1^t F_1^j] & \mathbb{E} [F_1^t F_2^j] & 0 & \cdots & \cdots & 0 \\ \mathbb{E} [F_2^t F_1^j] & \mathbb{E} [F_2^t F_2^j] & \mathbb{E} [F_2^t F_3^j] & \ddots & \ddots & \vdots \\ 0 & \mathbb{E} [F_3^t F_2^j] & \mathbb{E} [F_3^t F_3^j] & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \mathbb{E} [F_{r-1}^t F_{r-1}^j] & \mathbb{E} [F_{r-1}^t F_r^j] \\ 0 & \cdots & \cdots & 0 & \mathbb{E} [F_r^t F_{r-1}^j] & \mathbb{E} [F_r^t F_r^j] \end{bmatrix}. \end{aligned}$$

By law of large numbers,

$$\begin{aligned} \mathbf{V} &= \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_t, \beta)'] + \lim_{N \rightarrow \infty} \frac{N-1}{N} \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t+1}, \beta)'] \\ &\quad + \lim_{N \rightarrow \infty} \frac{N-1}{N} \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t-1}, \beta)'] \\ &= \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_t, \beta)'] + \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t+1}, \beta)'] + \mathbb{E} [F(\mathbf{X}_t, \beta) F(\mathbf{X}_{t-1}, \beta)'] \\ &= \mathbb{E} [G_t G_t'] + \mathbb{E} [G_t G_{t+1}'] + \mathbb{E} [G_t G_{t-1}'] \end{aligned}$$

and we can estimate  $\mathbf{V}$  by the sample counterpart.

Note that 2SLS attains the efficiency bound.

## 7.6 Testing with GMM

### 7.6.1 Calibration & Verification

1. **Calibrate** (= estimate) using  $k$  equations:

$$\mathbb{E} [F_1(\mathbf{X}_t, \beta)] = \mathbf{0}$$

where  $F_1$  has  $k$  coordinates.

2. **Verify** (= test) using  $r - k$  equations:

$$\mathbb{E} [F_2(\mathbf{X}_t, \boldsymbol{\beta})] = \mathbf{0}$$

while taking into account the initial estimation.

▷ To test this, use the fact that

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F_2(X_t, \mathbf{b}_N) \approx B' \left[ I - D (A'D)^{-1} A' \right] \frac{1}{\sqrt{N}} \sum_{t=1}^N F(X_t, \boldsymbol{\beta})$$

where  $A' = \begin{bmatrix} I & 0 \end{bmatrix}$  and  $B' = \begin{bmatrix} 0 & I \end{bmatrix}$ .

▷ Rewriting the above, partition  $D' = \begin{bmatrix} D_1 & D_2 \end{bmatrix}$  to write:

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F_2(X_t, \mathbf{b}_N) \approx \begin{bmatrix} -D_2 (D_1)^{-1} & I \end{bmatrix} \frac{1}{\sqrt{N}} \sum_{t=1}^N F(X_t, \boldsymbol{\beta})$$

and use a chi-square test.

### 7.6.2 Testing with Efficient Estimation

1. Set up the moment restriction and use the efficient selection matrix  $A' = D'V^{-1}$ .
2. Denote  $\Lambda$  such that  $V^{-1} = \Lambda' \Lambda$ . Since:

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F_2(X_t, \mathbf{b}_N) \approx B' \left[ I - D (A'D)^{-1} A' \right] \frac{1}{\sqrt{N}} \sum_{t=1}^N F(X_t, \boldsymbol{\beta})$$

Multiplying each side by  $\Lambda$  and plugging in  $A' = D'V^{-1}$ :

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F_2(X_t, \mathbf{b}_N) = B' \left[ I - \Lambda D (D' \Lambda' \Lambda D)^{-1} D' \Lambda' \right] \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(X_t, \boldsymbol{\beta})$$

Denote  $\Delta = \Lambda D (D' \Lambda' \Lambda D)^{-1} D' \Lambda'$  and recognize that  $\Delta$  and  $I - \Delta$  are both idempotent.  $\Delta$  has rank  $k$  and  $I - \Delta$  has rank  $r - k$ .

3. Then we have:

$$\left( \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) \right)' \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) \Rightarrow \chi^2(r)$$

which we can split into two components:

$$\begin{aligned} \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) \right)' (I - \Delta) \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) &\Rightarrow \chi^2(r - k) \\ \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) \right)' \Delta \frac{1}{\sqrt{N}} \sum_{t=1}^N \Lambda F(\mathbf{X}_t, \boldsymbol{\beta}) &\Rightarrow \chi^2(k) \end{aligned}$$

4. If you replace  $\boldsymbol{\beta}$  with  $\mathbf{b}_N$ , you lose degrees of freedom.

## 7.7 Arguing for Consistency

We want to show

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \rightarrow \mathbb{E}[F(\mathbf{X}_t, \mathbf{b})]$$

i.e. a functional version of LLN.

## 7.8 Application: Exchange Rates

We consider regressions of the form:

$$\log S_{t+\ell} = \beta_0 + \beta_1 \log F_t + U_{t+\ell}$$

with  $\mathbb{E}[U_{t+\ell} | \mathfrak{F}_t] = 0$ . It's not actually a regression since we don't have exogenous regressors i.e. regressors are pre-determined. It turns out that you can still do least squares since

$$\begin{aligned} \mathbb{E}[U_{t+\ell}] &= 0 \\ \mathbb{E}[\log F_t U_{t+\ell}] &= 0 \end{aligned}$$

and we can use GMM to work out the OLS distribution.

Can you do better? We know that combining the properties above, we have

$$\mathbb{E}[Z_t U_{t+\ell}] = 0, \forall Z_t \in \mathfrak{F}_t$$

so I can put in other stuff. The  $U$ s will be serially correlated i.e.  $U_{t+\ell}$  and  $U_{t+\ell-1}$  will be correlated. You can put in lag regressors and improve efficiency. What you cannot do is put regressors going forward ( $F_{t+1}, F_{t+2}, \dots$ ). GLS can let you go both directions.

## 7.9 Key GMM Equations

**Baseline restriction:**

$$\mathbb{E}[F(\mathbf{X}, \mathbf{b})] = \mathbf{0}_{r \times 1} \text{ if and only if } \mathbf{b} = \beta$$

**Sample counterpart of the moment condition**

$$\mathbf{A}' \frac{1}{N} \sum_{t=1}^N F_t(\mathbf{X}_t, \mathbf{b}_N) = \mathbf{0}.$$

**GMM asymptotic normality**

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) \sim N(0, \mathbf{V}),$$

where

$$\begin{aligned} \mathbf{V} &= \mathbb{E}[\mathbf{H}_{t+1} \mathbf{H}_{t+1}'] \\ &\approx \lim_{N \rightarrow \infty} \mathbb{E} \left[ \left( \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \beta) \right) \left( \frac{1}{\sqrt{N}} \sum_{j=1}^N F(\mathbf{X}_j, \beta) \right)' \right] \end{aligned}$$

and  $\mathbf{H}_{t+1}$  is the martingale component of  $F(\mathbf{X}_t, \beta)$ .

**Efficient selection matrix**

$$\mathbf{A}^* = \mathbf{V}^{-1} \mathbf{D},$$

where  $\mathbf{D}_{r \times k} = \mathbb{E} [\partial F(\mathbf{X}_t, \boldsymbol{\beta}) / \partial \mathbf{b}']$ . Note that  $\mathbf{V}$  also has the interpretation of the weighting matrix.

**Approximation**

$$\sqrt{N}(\mathbf{b}_N - \boldsymbol{\beta}) \approx -(\mathbf{A}'\mathbf{D})^{-1} \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta})$$

**Hypothesis testing with GMM** For both tests outlined below, you construct Wald tests that is distributed according to  $\chi^2$ .

1. Testing coefficient estimates:

$$\sqrt{N}(\mathbf{b}_N - \boldsymbol{\beta}) \approx -(\mathbf{A}'\mathbf{D})^{-1} \mathbf{A}' \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta})$$

2. Testing moment conditions:

$$\frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N) \approx \left( \mathbb{I}_r - \mathbf{D}(\mathbf{A}'\mathbf{D})^{-1} \mathbf{A}' \right) \frac{1}{\sqrt{N}} \sum_{t=1}^N F(\mathbf{X}_t, \boldsymbol{\beta})$$

Note that premultiplying the RHS matrix is equal to zero:

$$\mathbf{A}' \left( \mathbb{I}_r - \mathbf{D}(\mathbf{A}'\mathbf{D})^{-1} \mathbf{A}' \right) = \mathbf{0}$$

and this matrix is adjusting for the estimation of  $\boldsymbol{\beta}$ .

**7.10 Revisiting Linear Models via GMM****7.10.1 Base line**

Let  $\mathbf{Y}_t \in \mathbb{R}^{k+1}$ ,  $\boldsymbol{\beta} \in \mathbb{R}^k$ , and  $u_t \in \mathbb{R}$  such that

$$\mathbf{Y}_t' \begin{bmatrix} 1 \\ \boldsymbol{\beta} \end{bmatrix} = u_t.$$

Our goal is to estimate  $\boldsymbol{\beta}$  consistently. Let  $\mathfrak{F}_{t-\ell}$  denote the information set in period  $t - \ell$  and  $\mathbf{Z}_{t-\ell} \in \mathfrak{F}_{t-\ell} \subset \mathbb{R}^r$ .

**Moment condition** We suppose that, for some  $\ell$ ,

$$\mathbb{E}[u_t | \mathfrak{F}_\tau] \begin{cases} = 0 & \text{if } \tau \leq t - \ell, \\ \neq 0 & \text{if } t - \ell < \tau \leq t + \ell. \end{cases}$$

If  $\ell = 0$ , then

$$\begin{aligned} \mathbb{E}[u_t | \mathfrak{F}_t] &= 0, \\ \mathbb{E}[u_t | \mathfrak{F}_\tau] &\neq 0, \forall \tau \neq t. \end{aligned}$$

If, instead,  $\ell = 1$ , then

$$\begin{aligned} \mathbb{E}[u_t | \mathfrak{F}_{t-1}] &= \mathbb{E}[u_t | \mathfrak{F}_{t-2}] = \cdots = 0, \\ \mathbb{E}[u_t | \mathfrak{F}_t], \mathbb{E}[u_t | \mathfrak{F}_{t+1}] &\neq 0. \end{aligned}$$

Thus, the choice of  $\ell$  determines the correlation between the error term  $u_\tau$  for  $t - \ell < \tau \leq t + \ell$ . The above conditional moment restriction implies that

$$\begin{aligned} \mathbb{E}[u_t | \mathfrak{F}_{t-\ell}] &= 0 \\ \Rightarrow \mathbf{Z}_{t-\ell} \mathbb{E}[u_t | \mathfrak{F}_{t-\ell}] &= \mathbf{Z}_{t-\ell} 0 \\ \Leftrightarrow \mathbb{E}[\mathbf{Z}_{t-\ell} u_t | \mathfrak{F}_{t-\ell}] &= \mathbf{0}_{r \times 1} \\ \Rightarrow \mathbb{E}[\mathbb{E}[\mathbf{Z}_{t-\ell} u_t | \mathfrak{F}_{t-\ell}]] &= \mathbf{0}_{r \times 1} \\ \Leftrightarrow \mathbb{E}[\mathbf{Z}_{t-\ell} u_t] &= \mathbf{0}_{r \times 1}; \end{aligned}$$

i.e. it implies that  $\mathbf{Z}_{t-\ell}$  is orthogonal to  $u_t$ . Let

$$\mathbf{X}_t := \begin{bmatrix} \mathbf{Y}_t \\ \mathbf{Z}_{t-\ell} \end{bmatrix}$$

and partition  $\mathbf{Y}_t$  as

$$\mathbf{Y}_t = \begin{bmatrix} Y_t^1 \\ \mathbf{Y}_t^{2'} \end{bmatrix},$$

where  $Y_t^1$  is a scalar and  $\mathbf{Y}_t^{2'}$  is an  $k \times 1$  vector. Define

$$\begin{aligned} F(\mathbf{X}, \boldsymbol{\beta})_{r \times 1} &:= \mathbf{Z}_{t-\ell} u_t \\ &= \mathbf{Z}_{t-\ell} \mathbf{Y}_t' \begin{bmatrix} 1 \\ \boldsymbol{\beta} \end{bmatrix} \\ &= \mathbf{Z}_{t-\ell} \begin{bmatrix} Y_t^1 & \mathbf{Y}_t^{2'} \end{bmatrix} \begin{bmatrix} 1 \\ \boldsymbol{\beta} \end{bmatrix} \\ &= \mathbf{Z}_{t-\ell} (Y_t^1 + \mathbf{Y}_t^{2'} \boldsymbol{\beta}). \end{aligned}$$

We can then write the orthogonality condition, (??), as

$$\mathbb{E}[F(\mathbf{X}, \boldsymbol{\beta})] = \mathbf{0}_{r \times 1} = \mathbb{E}[\mathbf{Z}_{t-\ell} u_t].$$

Thus, we see that the GMM moment condition coincides with the orthogonality condition for IV. The sample analog of the moment restriction is

$$\begin{aligned} \mathbf{0}_{r \times 1} &= \frac{1}{N} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \\ &= \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}). \end{aligned}$$

**Derivative matrix** Consider

$$\frac{\partial F(\mathbf{X}_t, \boldsymbol{\beta})}{\partial \mathbf{b}'} = \frac{\partial \mathbf{Z}_{t-\ell} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b})}{\partial \mathbf{b}'} = \mathbf{Z}_{t-\ell} \mathbf{Y}_t^{2'}.$$

Since  $F$  is linear, we realise that the derivative does not involve  $\boldsymbol{\beta}$ . Then,

$$\mathbf{D} := \mathbb{E} \left[ \frac{\partial F(\mathbf{X}_t, \boldsymbol{\beta})}{\partial \mathbf{b}'} \right] = \mathbb{E} [\mathbf{Z}_{t-\ell} \mathbf{Y}_t^{2'}].$$

The sample analog,  $\mathbf{D}_N$ , is then

$$\mathbf{D}_N = \frac{1}{N} \sum_{t=1}^N \frac{\partial F(\mathbf{X}_t, \mathbf{b})}{\partial \mathbf{b}'} = \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} \mathbf{Y}_t^{2'}.$$

Note  $\mathbf{D}_N \xrightarrow{p} \mathbf{D}$ .



**Selection matrix** We have  $r$  linear equations and  $k$  unknowns. The selection matrix,  $\mathbf{A}$ , thus has dimension  $r \times k$ .

**Estimating  $\beta$**  Premultiplying the sample analog of the moment restriction by the (transpose of the) selection matrix yields

$$\begin{aligned} \mathbf{0}_{k \times 1} &= \mathbf{A}' \frac{1}{N} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}_N) \\ &= \mathbf{A}' \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}_N) \\ \Leftrightarrow -\mathbf{A}' \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} \mathbf{Y}_t^{2'} \mathbf{b}_N &= \mathbf{A}' \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} Y_t^1 \\ \Leftrightarrow \mathbf{b}_N &= - \left( \mathbf{A}' \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} \mathbf{Y}_t^{2'} \right)^{-1} \left( \mathbf{A}' \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-\ell} Y_t^1 \right). \end{aligned}$$

We can see that the estimator  $\mathbf{b}_N$  depends on the selection matrix  $\mathbf{A}$ .

### 7.10.2 Least squares

Given partition (??), we can write

$$\begin{aligned} \mathbf{Y}_t' \begin{bmatrix} 1 \\ \beta \end{bmatrix} &= u_t \Leftrightarrow Y_t^1 + \mathbf{Y}_t^{2'} \beta = u_t \\ &\Leftrightarrow Y_t^1 = -\mathbf{Y}_t^{2'} \beta + u_t, \end{aligned}$$

which is the familiar regression form. We can then impose that

$$\mathbb{E}[u_t | \mathfrak{F}_t] = 0,$$

which implies that

$$\mathbb{E}[\mathbf{Y}_t u_t] = 0.$$

The sample analog of the moment restriction is

$$\begin{aligned} \mathbf{0}_{r \times 1} &= \frac{1}{N} \sum_{t=1}^N F(\mathbf{X}_t, \mathbf{b}) \\ &= \frac{1}{N} \sum_{t=1}^N \mathbf{Y}_t (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}) \\ &= \frac{1}{N} \sum_{t=1}^N \begin{bmatrix} Y_t^1 \\ \mathbf{Y}_t^2 \end{bmatrix} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}). \end{aligned}$$

In this case,  $r = k + 1$ . Now, let

$$\mathbf{A} = \begin{bmatrix} 0 & \mathbf{0}_{1 \times k} \\ \mathbf{0}_{k \times 1} & \mathbb{I}_k \end{bmatrix}_{r \times (k+1)}.$$

This selection matrix selects only  $\mathbf{Y}_t^2$  from  $\mathbf{Y}_t$ . Then,

$$\begin{aligned}
 \mathbf{A}'\mathbf{0}_{r \times 1} &= \mathbf{A}' \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} Y_t^1 \\ \mathbf{Y}_t^2 \end{bmatrix} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}_N) \\
 \Leftrightarrow \mathbf{0}_{(k+1) \times 1} &= \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} 0 & \mathbf{0}_{1 \times k} \\ \mathbf{0}_{k \times 1} & \mathbb{I}_k \end{bmatrix} \begin{bmatrix} Y_t^1 \\ \mathbf{Y}_t^2 \end{bmatrix} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}_N) \\
 &= \frac{1}{N} \sum_{i=1}^N \begin{bmatrix} 0 \\ \mathbf{Y}_t^2 \end{bmatrix} (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}_N) \\
 \Rightarrow \mathbf{0}_{k \times 1} &= \frac{1}{N} \sum_{i=1}^N \mathbf{Y}_t^2 (Y_t^1 + \mathbf{Y}_t^{2'} \mathbf{b}_N),
 \end{aligned}$$

which is the usual OLS condition. Rearranging yields that

$$\mathbf{b}_N = \left( \frac{1}{N} \sum_{i=1}^N (-\mathbf{Y}_t^2) (-\mathbf{Y}_t^2)' \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N (-\mathbf{Y}_t^2) Y_t^1 \right).$$

### 7.10.3 Two-stage least squares

Suppose  $\ell = 1$  and that

$$\mathbb{E} [u_t^2 | \mathfrak{F}_{t-1}] = \sigma^2.$$

Thus, we are assuming no serial correlation and homoscedasticity. The moment restriction is, as before,

$$\mathbf{0}_{r \times 1} = \mathbb{E} [\mathbf{Z}_{t-1} u_t].$$

**Efficient selection matrix** Recall that the efficient selection matrix is given by

$$\mathbf{A}^* = \mathbf{V}^{-1} \mathbf{D}.$$

(we can also multiply this by any conforming nonsingular matrix.) As derived before, the derivative matrix is given by

$$\mathbf{D} := \mathbb{E} \left[ \frac{\partial F(\mathbf{X}_t, \boldsymbol{\beta})}{\partial \mathbf{b}'} \right] = \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Y}_t^{2'}].$$

To obtain  $\mathbf{V}$ , we construct an additive martingale  $\{H_\tau\}$  as

$$\begin{aligned}
 H_\tau &:= \sum_{t=1}^{\tau} \mathbf{Z}_{t-1} u_t = \sum_{t=1}^{\tau} \mathbf{Z}_{t-1} (Y_t^1 + \mathbf{Y}_t^{2'} \boldsymbol{\beta}) \\
 \Rightarrow H_\tau - H_{\tau-1} &= \mathbf{Z}_{\tau-1} u_\tau = \kappa(\mathbf{Z}_{\tau-1}, u_\tau),
 \end{aligned}$$

so that  $\mathbf{Z}_{\tau-1} u_\tau$  is the additive increment. Given the moment restriction, observe that

$$\mathbb{E} [\kappa(\mathbf{Z}_{\tau-1}, u_\tau) | \mathbf{Z}_{\tau-1}] = \mathbf{Z}_{\tau-1} \mathbb{E} [u_\tau] = \mathbf{0}_{r \times 1}$$

so that we verify that  $\{H_\tau\}$  is indeed an additive martingale. Then, by Billingsley Central Limit Theorem,

$$\sqrt{N} \left( \frac{1}{N} \sum_{t=1}^N \mathbf{Z}_{t-1} u_t \right) \xrightarrow{d} N(0, \mathbf{V}),$$

where

$$\begin{aligned}\mathbf{V} &= \mathbb{E} \left[ (\mathbf{Z}_{t-1} u_t) (\mathbf{Z}_{t-1} u_t)' \right] \\ &= \mathbb{E} \left[ \mathbf{Z}_{t-1} \mathbb{E} [u_t^2 | \mathfrak{F}_{t-1}] \mathbf{Z}_{t-1}' \right] \\ &= \sigma^2 \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Z}_{t-1}'] .\end{aligned}$$

Then,

$$\mathbf{A}^* = \frac{1}{\sigma^2} \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Z}_{t-1}']^{-1} \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Y}_t^{2'}] .$$

Recall that

$$\text{Cov} [\mathbf{A}^*] = \text{Cov} [\mathbf{A}^* \mathbf{B}]$$

for any nonsingular matrix  $\mathbf{B}$ . Letting

$$\mathbf{B} = \sigma^2 \mathbb{I},$$

then

$$\tilde{\mathbf{A}}^* := \mathbf{A}^* \mathbf{B} = \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Z}_{t-1}']^{-1} \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Y}_t^{2'}] .$$

This gives us an efficient selection matrix. Observe that this is coefficient on  $\mathbf{Z}_{t-1}'$  on the population regression of  $\mathbf{Y}_t^{2'}$  on  $\mathbf{Z}_{t-1}$ . That is, the selection matrix is coefficient from the first-stage regression in 2SLS.

**Relationship between GMM and IV estimators** Let us consider the population regression first. The moment restriction is

$$\mathbf{0}_{r \times 1} = \mathbb{E} [\mathbf{Z}_{t-1} (Y_t^1 + \mathbf{Y}_t^{2'} \beta)] .$$

Premultiplying by  $\tilde{\mathbf{A}}^{*'}$

$$\begin{aligned}\mathbf{0}_{k \times 1} &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} (Y_t^1 + \mathbf{Y}_t^{2'} \beta)] \\ \Leftrightarrow \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} \mathbf{Y}_t^{2'}] \beta &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1] .\end{aligned}$$

Consider the best linear predictor of  $\mathbf{Y}_t^{2'} | \mathbf{Z}_{t-1}$ :

$$\mathbf{Y}_t^2 = \mathbf{\Pi}' \mathbf{Z}_{t-1} + \mathbf{v}_t ,$$

where, as noted before,  $\mathbf{\Pi} = \tilde{\mathbf{A}}^*$ . By the property of BLP,

$$\mathbb{E} [\mathbf{Z}_{t-1} \mathbf{v}_t'] = \mathbf{0} .$$

This means that

$$\begin{aligned}\mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} (\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} + \mathbf{v}_t)'] \beta &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1] \\ \Leftrightarrow \left( \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}' \tilde{\mathbf{A}}^*] + \tilde{\mathbf{A}}^{*'} \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{v}_t'] \right) \beta &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1] \\ \Rightarrow \left( \tilde{\mathbf{A}}^{*'} \mathbb{E} [\mathbf{Z}_{t-1} \mathbf{Z}_{t-1}'] \tilde{\mathbf{A}}^* \right) \beta &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1] \\ \Leftrightarrow \beta &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}' \tilde{\mathbf{A}}^*]^{-1} \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1] \\ &= \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}' \tilde{\mathbf{A}}^*]^{-1} \mathbb{E} [\tilde{\mathbf{A}}^{*'} \mathbf{Z}_{t-1} Y_t^1]\end{aligned}$$

This gives us the familiar IV (population) estimator (over-identified case). Hence, in this case, the GMM estimator and the IV estimator coincide.

#### 7.10.4 Heteroscedasticity robust estimator

Now suppose that  $\mathbb{E}[u_t^2 | \mathfrak{F}_{t-1}]$  is not a constant. The only difference is that we can no longer simplify  $\mathbf{V}$  as we did before:

$$\begin{aligned}\mathbf{V} &= \mathbb{E}[(\mathbf{Z}_{t-1} u_t)(\mathbf{Z}_{t-1} u_t)'] \\ &= \mathbb{E}[u_t^2 \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}'] .\end{aligned}$$

Then,

$$\mathbf{A}^* = \mathbb{E}[u_t^2 \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}']^{-1} \mathbb{E}[\mathbf{Z}_{t-1} \mathbf{Y}_t^{2'}] .$$

We can estimate  $\mathbf{V}$  as

$$\hat{\mathbf{V}} = \frac{1}{N} \sum_{t=1}^N \hat{u}_t^2 \mathbf{Z}_{t-1} \mathbf{Z}_{t-1}' ,$$

where  $\hat{u}_t = Y_t^1 + \mathbf{Y}_t^{2'} \hat{\boldsymbol{\beta}}$ . The TSLS estimator is no longer efficient but remains consistent.