

Appendix A.

Functional Analysis

This appendix provides an introduction to the analysis of functional equations (functional analysis). It describes the contraction mapping theorem, a workhorse for studying dynamic programs.

A.1. Metric spaces and operators

We begin with the definition of a metric space, which is a pair of objects, a set X , and a function d .¹

Definition A.1.1. *A metric space is a set X and a function d called a metric, $d: X \times X \rightarrow \mathbb{R}$. The metric $d(x, y)$ satisfies the following four properties:*

- M1. *Positivity:* $d(x, y) \geq 0$ for all $x, y \in X$.
- M2. *Strict positivity:* $d(x, y) = 0$ if and only if $x = y$.
- M3. *Symmetry:* $d(x, y) = d(y, x)$ for all $x, y \in X$.
- M4. *Triangle inequality:* $d(x, y) \leq d(x, z) + d(z, y)$ for all x, y , and $z \in X$.

We give some examples of the metric spaces with which we will be working:

Example A.1. $l_p[0, \infty)$. We say that $X = l_p[0, \infty)$ is the set of all sequences of complex numbers $\{x_t\}_{t=0}^{\infty}$ for which $\sum_{t=0}^{\infty} |x_t|^p$ converges, where $1 \leq p < \infty$. The function $d_p(x, y) = (\sum_{t=0}^{\infty} |x_t - y_t|^p)^{1/p}$ is a metric. Often, we will say that $p = 2$ and will work in $l_2[0, \infty)$.

Example A.2. $l_{\infty}[0, \infty)$. The set $X = l_{\infty}[0, \infty)$ is the set of bounded sequences $\{x_t\}_{t=0}^{\infty}$ of real or complex numbers. The metric is $d_{\infty}(x, y) = \sup_t |x_t - y_t|$.

Example A.3. $l_p(-\infty, \infty)$ is the set of “two-sided” sequences $\{x_t\}_{t=-\infty}^{\infty}$ such that $\sum_{t=-\infty}^{\infty} |x_t|^p < +\infty$, where $1 \leq p < \infty$. The associated metric is $d_p(x, y) = (\sum_{t=-\infty}^{\infty} |x_t - y_t|^p)^{1/p}$.

¹ General references on the mathematics described in this appendix are Luenberger (1969) and Naylor and Sell (1982).

Example A.4. $l_\infty(-\infty, \infty)$ is the set of bounded sequences $\{x_t\}_{t=-\infty}^\infty$ with metric $d_\infty(x, y) = \sup |x_t - y_t|$.

Example A.5. Let $X = C[0, T]$ be the set of all continuous functions mapping the interval $[0, T]$ into R . We consider the metric

$$d_p(x, y) = \left[\int_0^T |x(t) - y(t)|^p dt \right]^{1/p},$$

where the integration is in the Riemann sense.

Example A.6. Let $X = C[0, T]$ be the set of all continuous functions mapping the interval $[0, T]$ into R . We consider the metric

$$d_\infty(x, y) = \sup_{0 \leq t \leq T} |x(t) - y(t)|.$$

We now have the following important definition:

Definition A.1.2. A sequence $\{x_n\}$ in a metric space (X, d) is said to be a Cauchy sequence if for each $\epsilon > 0$ there exists an $N(\epsilon)$ such that $d(x_n, x_m) < \epsilon$ for any $n, m \geq N(\epsilon)$. Thus a sequence $\{x_n\}$ is said to be Cauchy if $\lim_{n, m \rightarrow \infty} d(x_n, x_m) = 0$.

We also have the following definition of convergence:

Definition A.1.3. A sequence $\{x_n\}$ in a metric space (X, d) is said to converge to a limit $x_0 \in X$ if for every $\epsilon > 0$ there exists an $N(\epsilon)$ such that $d(x_n, x_0) < \epsilon$ for $n \geq N(\epsilon)$.

The following lemma asserts that every convergent sequence in (X, d) is a Cauchy sequence:

Lemma A.1.1. Let $\{x_n\}$ be a convergent sequence in a metric space (X, d) . Then $\{x_n\}$ is a Cauchy sequence.

Proof. Fix any $\epsilon > 0$. Let $x_0 \in X$ be the limit of $\{x_n\}$. Then for all m, n one has

$$d(x_n, x_m) \leq d(x_n, x_0) + d(x_m, x_0)$$

by virtue of the triangle inequality. Because x_0 is the limit of $\{x_n\}$, there exists an N such that $d(x_n, x_0) < \epsilon/2$ for $n \geq N$. Together with the preceding

inequality, this statement implies that $d(x_n, x_m) < \epsilon$ for $n, m \geq N$. Therefore, $\{x_n\}$ is a Cauchy sequence. ■

We now consider two examples of sequences in metric spaces. The examples are designed to illustrate aspects of the concept of a Cauchy sequence. We first consider the metric space $\{C[0, 1], d_2(x, y)\}$. We let $\{x_n\}$ be the sequence of continuous functions $x_n(t) = t^n$. Evidently this sequence converges pointwise to the function

$$x_0(t) = \begin{cases} 0, & 0 \leq t < 1 \\ 1, & t = 1. \end{cases}$$

Now, in $\{C[0, 1], d_2(x, y)\}$, the sequence $x_n(t)$ is a Cauchy sequence. To verify this claim, calculate

$$d_2(t^m, t^n)^2 = \int_0^1 (t^n - t^m)^2 dt = \frac{1}{2n+1} + \frac{1}{2m+1} - \frac{2}{m+n+1}.$$

Clearly, for any $\epsilon > 0$, it is possible to choose an $N(\epsilon)$ that makes the square root of the right side less than ϵ whenever m and n both exceed N . Thus $x_n(t)$ is a Cauchy sequence. Notice, however, that the limit point $x_0(t)$ does not belong to $\{C[0, 1], d_2(x, y)\}$ because it is not a continuous function.

As our second example, we consider the space $\{C[0, 1], d_\infty(x, y)\}$. We consider the sequence $x_n(t) = t^n$. In $(C[0, 1], d_\infty)$, the sequence $x_n(t)$ is not a Cauchy sequence. To verify this claim, it is sufficient to establish that, for any fixed $m > 0$, there is a $\delta > 0$ such that

$$\sup_{n>0} \sup_{0 \leq t \leq 1} |t^n - t^m| > \delta.$$

Direct calculations show that, for fixed m ,

$$\sup_n \sup_{0 \leq t \leq 1} |t^n - t^m| = 1.$$

Parenthetically we may note that

$$\begin{aligned} \sup_{n>0} \sup_{0 \leq t \leq 1} |t^n - t^m| &= \sup_{0 \leq t \leq 1} \sup_{n>0} |t^n - t^m| = \sup_{0 \leq t \leq 1} \lim_{n \rightarrow \infty} |t^n - t^m| \\ &= \sup_{0 \leq t \leq 1} \lim_{n \rightarrow \infty} t^m |t^{n-m} - 1| = \sup_{0 \leq t \leq 1} t^m = 1. \end{aligned}$$

Therefore, $\{t^n\}$ is not a Cauchy sequence in $(C[0, 1], d_\infty)$.

These examples illustrate the fact that whether a given sequence is Cauchy depends on the metric space within which it is embedded, in particular on the metric that is being used. The sequence $\{t^n\}$ is Cauchy in $(C[0, 1], d_2)$, and more generally in $(C[0, 1], d_p)$ for $1 \leq p < \infty$. The sequence $\{t^n\}$, however, is *not* Cauchy in the metric space $(C[0, 1], d_\infty)$. The first example also illustrates the fact that a Cauchy sequence in (X, d) need *not* converge to a limit point x_0 belonging to the metric space. The property that Cauchy sequences converge to points lying in the metric space is desirable in many applications. We give this property a name.

Definition A.1.4. *A metric space (X, d) is said to be complete if each Cauchy sequence in (X, d) is a convergent sequence in (X, d) . That is, in a complete metric space, each Cauchy sequence converges to a point belonging to the metric space.*

The following metric spaces are complete:

$$\begin{aligned} (l_p[0, \infty), d_p), \quad & 1 \leq p < \infty \\ (l_\infty[0, \infty), d_\infty) \\ (C[0, T], d_\infty). \end{aligned}$$

The following metric spaces are not complete:

$$(C[0, T], d_p), \quad 1 \leq p < \infty.$$

Proofs that $(l_p[0, \infty), d_p)$ for $1 \leq p \leq \infty$ and $(C[0, T], d_\infty)$ are complete are contained in Naylor and Sell (1982, chap. 3). In effect, we have already shown by counterexample that the space $(C[0, 1], d_2)$ is not complete, because we displayed a Cauchy sequence that did not converge to a point in the metric space. A definition may now be stated:

Definition A.1.5. *A function f mapping a metric space (X, d) into itself is called an operator.*

We need a notion of continuity of an operator.

Definition A.1.6. *Let $f : X \rightarrow X$ be an operator on a metric space (X, d) . The operator f is said to be continuous at a point $x_0 \in X$ if for every $\epsilon > 0$*

there exists a $\delta > 0$ such that $d[f(x), f(x_0)] < \epsilon$ whenever $d(x, x_0) < \delta$. The operator f is said to be continuous if it is continuous at each point $x \in X$.

We shall be studying an operator with a particular property, the application of which to any two distinct points $x, y \in X$ brings them closer together.

Definition A.1.7. Let (X, d) be a metric space and let $f : X \rightarrow X$. We say that f is a contraction or contraction mapping if there is a real number $k, 0 \leq k < 1$, such that

$$d[f(x), f(y)] \leq kd(x, y) \quad \text{for all } x, y \in X.$$

It follows directly from the definition that a contraction mapping is a continuous operator.

We now state the following theorem:

Theorem A.1.1. *Contraction Mapping*

Let (X, d) be a complete metric space and let $f : X \rightarrow X$ be a contraction. Then there is a unique point $x_0 \in X$ such that $f(x_0) = x_0$. Furthermore, if x is any point in X and $\{x_n\}$ is defined inductively according to $x_1 = f(x), x_2 = f(x_1), \dots, x_{n+1} = f(x_n)$, then $\{x_n\}$ converges to x_0 .

Proof. Let x be any point in X . Define $x_1 = f(x), x_2 = f(x_1), \dots$. Express this as $x_n = f^n(x)$. To show that the sequence x_n is Cauchy, first assume that $n > m$. Then

$$\begin{aligned} d(x_n, x_m) &= d[f^n(x), f^m(x)] = d[f^m(x_{n-m}), f^m(x)] \\ &\leq kd[f^{m-1}(x_{n-m}), f^{m-1}(x)] \end{aligned}$$

By induction, we get

$$(*) \quad d(x_n, x_m) \leq k^m d(x_{n-m}, x).$$

When we repeatedly use the triangle inequality, the preceding inequality implies that

$$d(x_n, x_m) \leq k^m [d(x_{n-m}, x_{n-m-1}) + \dots + d(x_2, x_1) + d(x_1, x)].$$

Applying $(*)$ gives

$$d(x_n, x_m) \leq k^m (k^{n-m-1} + \dots + k + 1) d(x_1, x).$$

Because $0 \leq k < 1$, we have

$$(\dagger) \quad d(x_n, x_m) \leq k^m \sum_{i=0}^{\infty} k^i d(x_1, x) = \frac{k^m}{1-k} d(x_1, x).$$

The right side of (\dagger) can be made arbitrarily small by choosing m sufficiently large. Therefore, $d(x_n, x_m) \rightarrow 0$ as $n, m \rightarrow \infty$. Thus $\{x_n\}$ is a Cauchy sequence. Because (X, d) is complete, $\{x_n\}$ converges to an element of (X, d) .

The limit point x_0 of $\{x_n\} = \{f^n(x)\}$ is a fixed point of f . Because f is continuous, $\lim_{n \rightarrow \infty} f(x_n) = f(\lim_{n \rightarrow \infty} x_n)$. Now $f(\lim_{n \rightarrow \infty} x_n) = f(x_0)$ and $\lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} x_{n+1} = x_0$. Therefore $x_0 = f(x_0)$.

To show that the fixed point x_0 is unique, assume the contrary. Assume that x_0 and y_0 , $x_0 \neq y_0$, are two fixed points of f . But then

$$0 < d(x_0, y_0) = d[f(x_0), f(y_0)] \leq kd(x_0, y_0) < d(x_0, y_0),$$

which is a contradiction. Therefore f has a unique fixed point. ■

We now restrict ourselves to sets X whose elements are functions. The spaces $C[0, T]$ and $l_p[0, \infty)$ for $1 \leq p \leq \infty$ are examples of spaces of functions. Let us define the notion of inequality of two functions.

Definition A.1.8. Let X be a space of functions, and let $x, y \in X$. Then $x \geq y$ if and only if $x(t) \geq y(t)$ for every t in the domain of the functions.

Let X be a space of functions. We use the d_∞ metric, defined as $d_\infty(x, y) = \sup_t |x(t) - y(t)|$, where the supremum is over the domain of definition of the function.

A pair of conditions that are sufficient for an operator $T : (X, d_\infty) \rightarrow (X, d_\infty)$ to be a contraction appear in the following theorem:²

Theorem A.1.2. Blackwell's Sufficient Conditions for T to be a Contraction

Let T be an operator on a metric space (X, d_∞) , where X is a space of functions. Assume that T has the following two properties:

- (a) *Monotonicity:* For any $x, y \in X$, $x \geq y$ implies $T(x) \geq T(y)$.
- (b) *Discounting:* Let c denote a function that is constant at the real value c for

² See Blackwell's (1965) Theorem 5. This theorem is used extensively by Stokey and Lucas with Prescott (1989).

all points in the domain of definition of the functions in X . For any positive real c and every $x \in X$, $T(x+c) \leq T(x) + \beta c$ for some β satisfying $0 \leq \beta < 1$. Then T is a contraction mapping with modulus β .

Proof. For all $x, y \in X$, $x \leq y + d(x, y)$. Applying properties (a) and (b) to this inequality gives

$$T(x) \leq T(y + d(x, y)) \leq T(y) + \beta d(x, y).$$

Exchanging the roles of x and y and using the same logic implies

$$T(y) \leq T(x) + \beta d(x, y).$$

Combining these two inequalities gives $|T(x) - T(y)| \leq \beta d(x, y)$ or

$$d(T(x), T(y)) \leq \beta d(x, y).$$

■

A.2. Discounted dynamic programming

We study the functional equation associated with a discounted dynamic programming problem:

$$v(x) = \max_{u \in R^k} \{r(x, u) + \beta v(x')\}, \quad x' \leq g(x, u), \quad 0 < \beta < 1. \quad (A.2.1)$$

We assume that $r(x, u)$ is real valued, continuous, concave, and bounded and that the set $[x', x, u : x' \leq g(x, u), u \in R^k]$ is convex and compact.

We define the operator

$$Tv = \max_{u \in R^k} \{r(x, u) + \beta v(x')\}, \quad x' \leq g(x, u), \quad x \in X.$$

We work with the space of continuous bounded functions mapping X into the real line. We use the metric $d_\infty(v, w) = \sup_{x \in X} |v(x) - w(x)|$. This metric space is complete.

The operator T maps a continuous bounded function v into a continuous bounded function Tv . (For a proof, see Stokey and Lucas with Prescott, 1989.)³

We now establish that T is a contraction by verifying Blackwell's pair of sufficient conditions. First, suppose that $v(x) \geq w(x)$ for all $x \in X$. Then

$$\begin{aligned} Tv &= \max_{u \in R^k} \{r(x, u) + \beta v(x')\}, & x' \leq g(x, u) \\ &\geq \max_{u \in R^k} \{r(x, u) + \beta w(x')\}, & x' \leq g(x, u) \\ &= Tw. \end{aligned}$$

Thus, T is monotone. Next, notice that for any positive constant c ,

$$\begin{aligned} T(v + c) &= \max_{u \in R^k} \{r(x, u) + \beta[v(x') + c]\}, & x' \leq g(x, u) \\ &= \max_{u \in R^k} \{r(x, u) + \beta v(x') + \beta c\}, & x' \leq g(x, u) \\ &= Tv + \beta c. \end{aligned}$$

Thus, T discounts. Therefore, T satisfies both of Blackwell's conditions. It follows that T is a contraction on a complete metric space. Therefore the functional equation (A.2.1), which can be expressed as $v = Tv$, has a unique fixed point in the space of bounded continuous functions. This fixed point is approached in the limit in the d_∞ metric by iterations $v^n = T^n(v^0)$ starting from any bounded and continuous v^0 . Convergence in the d_∞ metric implies uniform convergence of the functions v^n .

Stokey and Lucas with Prescott (1989) show that T maps concave functions into concave functions. It follows that the solution of $v = Tv$ is a concave function.

³ The assertions in the preceding two paragraphs are the most difficult pieces of the argument to prove.

A.2.1. Policy improvement algorithm

For ease of exposition, in this section we shall assume that the constraint $x' \leq g(x, u)$ holds with equality. For the purposes of describing an alternative way to solve dynamic programming problems, we introduce a new operator. We use one step of iterating on the Bellman equation to define the new operator T_μ as follows:

$$T_\mu(v) = T(v)$$

or

$$T_\mu(v) = r[x, \mu(x)] + \beta v\{g[x, \mu(x)]\} ,$$

where $\mu(x)$ is the policy function that attains $T(v)(x)$. For a fixed $\mu(x)$, T_μ is an operator that maps bounded continuous functions into bounded continuous functions. Denote by C the space of bounded continuous functions mapping X into X .

For any admissible policy function $\mu(x)$, the operator T_μ is a contraction mapping. This fact can be established by verifying Blackwell's pair of sufficient conditions:

1. T_μ is monotone. Suppose that $v(x) \geq w(x)$. Then

$$\begin{aligned} T_\mu v &= r[x, \mu(x)] + \beta v\{g[x, \mu(x)]\} \\ &\geq r[x, \mu(x)] + \beta w\{g[x, \mu(x)]\} = T_\mu w . \end{aligned}$$

2. T_μ discounts. For any positive constant c

$$\begin{aligned} T_\mu(v + c) &= r(x, \mu) + \beta (v\{g[x, \mu(x)]\} + c) \\ &= T_\mu v + \beta c . \end{aligned}$$

Because T_μ is a contraction operator, the functional equation

$$v_\mu(x) = T_\mu[v_\mu(x)]$$

has a unique solution in the space of bounded continuous functions. This solution can be computed as a limit of iterations on T_μ starting from any bounded continuous function $v_0(x) \in C$,

$$v_\mu(x) = \lim_{k \rightarrow \infty} T_\mu^k(v_0)(x) .$$

The function $v_\mu(x)$ is the value of the objective function that would be attained by using the stationary policy $\mu(x)$ each period.

The following proposition describes the *policy iteration* or *Howard improvement* algorithm.

Theorem A.2.1. Let $v_\mu(x) = T_\mu[v_\mu(x)]$. Define a new policy $\bar{\mu}$ and an associated operator $T_{\bar{\mu}}$ by

$$T_{\bar{\mu}}[v_\mu(x)] = T[v_\mu(x)] ;$$

that is, $\bar{\mu}$ is the policy that solves a one-period problem with $v_\mu(x)$ as the terminal value function. Compute the fixed point

$$v_{\bar{\mu}}(x) = T_{\bar{\mu}}[v_{\bar{\mu}}(x)] .$$

Then $v_{\bar{\mu}}(x) \geq v_\mu(x)$. If $\mu(x)$ is not optimal, then $v_{\bar{\mu}}(x) > v_\mu(x)$ for at least one $x \in X$.

Proof. From the definition of $\bar{\mu}$ and $T_{\bar{\mu}}$, we have

$$\begin{aligned} T_{\bar{\mu}}[v_\mu(x)] &= r[x, \bar{\mu}(x)] + \beta v_\mu\{g[x, \bar{\mu}(x)]\} = \\ T(v_\mu)(x) &\geq r[x, \mu(x)] + \beta v_\mu\{g[x, \mu(x)]\} \\ &= T_\mu[v_\mu(x)] = v_\mu(x) \end{aligned}$$

or

$$T_{\bar{\mu}}[v_\mu(x)] \geq v_\mu(x) .$$

Apply $T_{\bar{\mu}}$ repeatedly to this inequality and use the monotonicity of $T_{\bar{\mu}}$ to conclude

$$v_{\bar{\mu}}(x) = \lim_{n \rightarrow \infty} T_{\bar{\mu}}^n[v_\mu(x)] \geq v_\mu(x) .$$

This establishes the asserted inequality $v_{\bar{\mu}}(x) \geq v_\mu(x)$. If $v_{\bar{\mu}}(x) = v_\mu(x)$ for all $x \in X$, then

$$\begin{aligned} v_\mu(x) &= T_{\bar{\mu}}[v_\mu(x)] \\ &= T[v_\mu(x)] , \end{aligned}$$

where the first equality follows because $T_{\bar{\mu}}[v_{\bar{\mu}}(x)] = v_{\bar{\mu}}(x)$, and the second equality follows from the definitions of $T_{\bar{\mu}}$ and $\bar{\mu}$. Because $v_\mu(x) = T[v_\mu(x)]$, the Bellman equation is satisfied by $v_\mu(x)$. ■

The *policy improvement* algorithm starts from an arbitrary feasible policy and iterates to convergence on the two following steps:⁴

⁴ A policy $\mu(x)$ is said to be *unimprovable* if it is optimal to follow it for the first period, given a terminal value function $v(x)$. In effect, the policy improvement algorithm starts with an arbitrary value function, then by solving a one-period problem, it generates an improved policy and an improved value function. The proposition states that optimality is characterized by the features, first, that there is no incentive to deviate from the policy during the first period, and second, that the terminal value function is the one associated with continuing the policy.

Step 1. For a feasible policy $\mu(x)$, compute $v_\mu = T_\mu(v_\mu)$.

Step 2. Find $\bar{\mu}$ by computing $T(v_\mu)$. Use $\bar{\mu}$ as the policy in step 1.

In many applications, this algorithm proves to be much faster than iterating on the Bellman equation.

A.2.2. A search problem

We now study the functional equation associated with a search problem of chapter 6. The functional equation is

$$v(w) = \max \left\{ \frac{w}{1-\beta}, \beta \int v(w') dF(w') \right\}, \quad 0 < \beta < 1. \quad (A.2.2)$$

Here, the wage offer drawn at t is w_t . Successive offers w_t are independently and identically distributed random variables. We assume that w_t has cumulative distribution function $\text{prob}\{w_t \leq w\} = F(w)$, where $F(0) = 0$ and $F(\bar{w}) = 1$ for some $\bar{w} < \infty$. In equation (A.2.2), $v(w)$ is the optimal value function for a currently unemployed worker who has offer w in hand. We seek a solution of the functional equation (A.2.2).

We work in the space of bounded continuous functions $C[0, \bar{w}]$ and use the d_∞ metric

$$d_\infty(x, y) = \sup_{0 \leq w \leq \bar{w}} |x(w) - y(w)|.$$

The metric space $(C[0, \bar{w}], d_\infty)$ is complete.

We consider the operator

$$T(z) = \max \left\{ \frac{w}{1-\beta}, \beta \int z(w') dF(w') \right\}. \quad (A.2.3)$$

Evidently the operator T maps functions z in $C[0, \bar{w}]$ into functions $T(z)$ in $C[0, \bar{w}]$. We now assert that the operator T defined by equation (A.2.3) is a contraction. To prove this assertion, we verify Blackwell's sufficient conditions. First, assume that $f(w) \geq g(w)$ for all $w \in [0, \bar{w}]$. Then note that

$$\begin{aligned} Tg &= \max \left\{ \frac{w}{1-\beta}, \beta \int g(w') dF(w') \right\} \\ &\leq \max \left\{ \frac{w}{1-\beta}, \beta \int f(w') dF(w') \right\} \\ &= Tf. \end{aligned}$$

Thus, T is monotone. Next, note that for any positive constant c ,

$$\begin{aligned}
 T(f + c) &= \max \left\{ \frac{w}{1 - \beta}, \beta \int [f(w') + c] dF(w') \right\} \\
 &= \max \left\{ \frac{w}{1 - \beta}, \beta \int f(w') dF(w') + \beta c \right\} \\
 &\leq \max \left\{ \frac{w}{1 - \beta}, \beta \int f(w') dF(w') \right\} + \beta c \\
 &= Tf + \beta c.
 \end{aligned}$$

Thus, T satisfies the discounting property and is therefore a contraction.

Application of the contraction mapping theorem, then, establishes that the functional equation $Tv = v$ has a unique solution in $C[0, \bar{w}]$, which is approached in the limit as $n \rightarrow \infty$ by $T^n(v^0) = v^n$, where v^0 is any point in $C[0, \bar{w}]$. Because the convergence in the space $C[0, \bar{w}]$ is in terms of the metric d_∞ , the convergence is uniform.

Appendix B.

Linear Projections and Hidden Markov Models

B.1. Linear projections

For reference we state the following theorems about linear least-squares projections. We let Y be an $(n \times 1)$ vector of random variables and X be an $(h \times 1)$ vector of random variables. We assume that the following first and second moments exist:

$$\begin{aligned} EY &= \mu_Y, \quad EX = \mu_X, \\ EXX' &= S_{XX}, \quad EYY' = S_{YY}, \quad EYX' = S_{YX}. \end{aligned}$$

Letting $x = X - EX$ and $y = Y - EY$, we define the following covariance matrices

$$Exx' = \Sigma_{xx}, \quad E'_{yy} = \Sigma_{yy}, \quad Eyx' = \Sigma_{yx}.$$

We are concerned with estimating Y as a linear function of X . The estimator of Y that is a linear function of X and that minimizes the mean squared error between each component Y and its estimate is called the *linear projection of Y on X* .

Definition B.1.1. The *linear projection* of Y on X is the affine function $\hat{Y} = AX + a_0$ that minimizes $E \text{ trace } \{(Y - \hat{Y})(Y - \hat{Y})'\}$ over all affine functions $a_0 + AX$ of X . We denote this linear projection as $\hat{E}[Y | X]$, or sometimes as $\hat{E}[Y | x, 1]$ to emphasize that a constant is included in the “information set.”

The linear projection of Y on X , $\hat{E}[Y | X]$ is also sometimes called the *wide sense expectation of Y conditional on X* . We have the following theorems:

Theorem B.1.1.

$$\hat{E}[Y | X] = \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(X - \mu_x). \quad (B.1.1)$$

Proof. The theorem follows immediately by writing out $E \text{ trace } (Y - \hat{Y})(Y - \hat{Y})'$ and completing the square, or else by writing out $E \text{ trace } (Y - \hat{Y})(Y - \hat{Y})'$ and obtaining first-order necessary conditions (“normal equations”) and solving them. ■

Theorem B.1.2.

$$\hat{E} \left[(Y - \hat{E}[Y | x]) | X' \right] = 0.$$

This equation states that the errors from the projection are orthogonal to each variable included in X .

Proof. Immediate from the normal equations. ■

Theorem B.1.3. (Orthogonality principle)

$$E \left[[Y - \hat{E}(Y | x)] x' \right] = 0.$$

Proof. Follows from Theorem 21.3. ■

Theorem B.1.4. (Orthogonal regressors)

Suppose that

$X' = (X_1, X_2, \dots, X_h)'$, $EX' = \mu' = (\mu_{x1}, \dots, \mu_{xh})'$, and $E(X_i - \mu_{xi})(X_j - \mu_{xj}) = 0$ for $i \neq j$. Then

$$\hat{E}[Y | x_1, \dots, x_n, 1] = \hat{E}[Y | x_1] + \hat{E}[Y | x_2] + \dots + \hat{E}[Y | x_n] - (n-1)\mu_y. \quad (B.1.2)$$

Proof. Note that from the hypothesis of orthogonal regressors, the matrix Σ_{xx} is diagonal. Applying equation (B.1.1) then gives equation (B.1.2). ■

B.2. Hidden Markov models

This section gives a brief introduction to hidden Markov models, a tool that is useful to study a variety of nonlinear filtering problems in finance and economics. We display a solution to a nonlinear filtering problem that a reader might want to compare to the linear filtering problem described earlier.

Consider an N -state Markov chain. We can represent the state space in terms of the unit vectors $S_x = \{e_1, \dots, e_N\}$, where e_i is the i th N -dimensional unit vector. Let the $N \times N$ transition matrix be P , with (i, j) element

$$P_{ij} = \text{Prob}(x_{t+1} = e_j \mid x_t = e_i).$$

With these definitions, we have

$$E x_{t+1} \mid x_t = P' x_t.$$

Define the “residual”

$$v_{t+1} = x_{t+1} - P' x_t,$$

which implies the linear “state-space” representation

$$x_{t+1} = P' x_t + v_{t+1}.$$

Notice how it follows that $E v_{t+1} \mid x_t = 0$, which qualifies v_{t+1} as a “martingale difference process adapted to x_t .”

We want to append a “measurement equation.” Suppose that x_t is not observed, but that y_t , a noisy function of x_t , is observed. Assume that y_t lives in the M -dimensional space S_y , which we represent in terms of M unit vectors: $S_y = \{f_1, \dots, f_M\}$, where f_i is the i th M -dimensional unit vector. To specify a linear measurement equation $y_t = C(x_t, u_t)$, where u_t is a measurement noise, we begin by defining the $N \times M$ matrix Q with

$$\text{Prob}(y_t = f_j \mid x_t = e_i) = Q_{ij}.$$

It follows that

$$E(y_t \mid x_t) = Q' x_t.$$

Define the residual

$$u_t \equiv y_t - E y_t \mid x_t,$$

which suggests the “observer equation”

$$y_t = Q'x_t + u_t.$$

It follows from the definition of u_t that $E u_t | x_t = 0$. Thus, we have the linear state-space system

$$\begin{aligned} x_{t+1} &= P'x_t + v_{t+1} \\ y_t &= Q'x_t + u_t. \end{aligned}$$

Using the definitions, it is straightforward to calculate the conditional second moments of the error processes v_{t+1}, u_t .¹

B.3. Nonlinear filtering

We seek a recursive formula for computing the conditional distribution of the hidden state:

$$\rho_i(t) = \text{Prob}\{x_t = i | y_1 = \eta_1, \dots, y_t = \eta_t\}.$$

Denote the history of observed y_t 's up to t as $\eta^t = \text{col}(\eta_1, \dots, \eta_t)$. Define the conditional probabilities

$$p(\xi_t, \eta_1, \dots, \eta_t) = \text{Prob}(x_t = \xi_t, y_1 = \eta_1, \dots, y_t = \eta_t),$$

¹ Notice that

$$\begin{aligned} x_{t+1}x'_{t+1} &= P'x_t(P'x_t)' + P'x_tv'_{t+1} \\ &\quad + v_{t+1}(P'x_t)' + v_{t+1}v'_{t+1} \end{aligned}$$

Substituting into this equation the facts that $x_{t+1}x'_{t+1} = \text{diag } x_{t+1} = \text{diag}(P'x_t) + \text{diag } v_{t+1}$ gives

$$\begin{aligned} v_{t+1}v'_{t+1} &= \text{diag}(P'x_t) + \text{diag}(v_{t+1}) - P'\text{diag } x_t P \\ &\quad - P'x_tv'_{t+1}(P'x_t)'. \end{aligned}$$

It follows that

$$E[v_{t+1}v'_{t+1} | x_t] = \text{diag}(P'x_t) - P'\text{diag } x_t P.$$

Similarly,

$$E[u_t u'_t | x_t] = \text{diag}(Q'x_t) - Q'\text{diag } x_t Q.$$

and assume $p(\eta_1, \dots, \eta_t) \neq 0$. Then apply the calculus of conditional expectations to compute²

$$\begin{aligned} p(\xi_t \mid \eta^t) &= \frac{p(\xi_t, \eta_t \mid \eta^{t-1})}{p(\eta_t \mid \eta^{t-1})} \\ &= \frac{\sum_{\xi_{t-1}} p(\eta_t \mid \xi_t) p(\xi_t \mid \xi_{t-1}) p(\xi_{t-1} \mid \eta^{t-1})}{\sum_{\xi_t} \sum_{\xi_{t-1}} p(\eta_t \mid \xi_t) p(\xi_t \mid \xi_{t-1}) p(\xi_{t-1} \mid \eta^{t-1})}. \end{aligned}$$

This result can be written

$$\rho_i(t+1) = \frac{\sum_s Q_{ij} P_{si} \rho_s(t)}{\sum_s \sum_i Q_{ij} P_{si} \rho_s(t)},$$

where $\eta_{t+1} = j$ is the value of y at $t+1$. We can represent this recursively as

$$\begin{aligned} \tilde{\rho}(t+1) &= \text{diag}(Q_j) P' \rho(t) \\ \rho(t+1) &= \frac{\tilde{\rho}(t+1)}{\langle \tilde{\rho}(t+1), \underline{1} \rangle}. \end{aligned}$$

where Q_j is the j th column of Q , and $\text{diag}(Q_j)$ is a diagonal matrix with Q_{ij} as the i th diagonal element; here $\langle \cdot, \cdot \rangle$ denotes the inner product of two vectors, and $\underline{1}$ is the unit vector.

² Notice that

$$\begin{aligned} p(\xi_t, \eta_t \mid \eta^{t-1}) &= \sum_{\xi_{t-1}} p(\xi_t, \eta_t, \xi_{t-1} \mid \eta^{t-1}) \\ &= \sum_{\xi_{t-1}} p(\xi_t, \eta_t \mid \xi_{t-1}, \eta^{t-1}) p(\xi_{t-1} \mid \eta^{t-1}) \\ p(\xi_t, \eta_t \mid \xi_{t-1}, \eta^{t-1}) &= p(\xi_t \mid \xi_{t-1}, \eta^{t-1}) p(\eta_t \mid \xi_t, \xi_{t-1}, \eta^{t-1}) \\ &= p(\xi_t \mid \xi_{t-1}) p(\eta_t \mid \xi_t). \end{aligned}$$

Combining these results gives the formula in the text.