

# Discrete-Time Stochastic Dynamic Programming

©1995, 1996 by Christopher Sims. This material may be freely reproduced for educational and research purposes, so long as it is not altered, this copyright notice is reproduced with it, and the copies are not sold.

## I. Notation and basic assumptions

We consider a problem defined in terms of

$t$  : a time index, with integer values

$C$  : a  $k \times 1$  vector, called the control vector

$S$  : an  $n \times 1$  vector, called the state vector

$\Gamma(\cdot)$  : a mapping from state vector values to subsets of  $R^k$ , defining constraints on the choice of  $C$

$\mathcal{I}_t$  : the information set at  $t$ , consisting of  $\{C(s), S(s), e(s), \text{all } s \leq t\}$

$e_t$  : a  $p \times 1$  random vector of disturbances at time  $t$ .

The objective is to maximize

$$E \left[ \sum_{t=0}^{\infty} b^t U(C_t, S_t) \right] \quad (1)$$

by choice of  $\{C_t, S_{t+1}, t = 0, \dots, \infty\}$ . We assume that the infinite sum inside the brackets in (1) is well-defined for each choice of  $C$ 's that satisfies the constraints below and that the expectation of the sum is well-defined for each such choice of  $C$ 's. The choice of  $C$ 's is constrained in four ways:

A)  $S_0$  is given, not subject to choice;

B) for each  $t=1, \dots, \infty$ ,  $S_t$  is determined from past history and current  $e_t$  according to

$$S_t = f(C_{t-1}, S_{t-1}, e_t) . \quad (2)$$

C) for each  $t$ ,  $C_t$  is constrained to lie in the set  $\Gamma(S_t)$ ;

D) for each  $t$ ,  $C_t$  is allowed to depend only on information in  $\mathcal{I}_t$ , and only in such a way that  $C_t$  and  $U(C_t, S_t)$  are well-defined random variables.

Note that (2) means that, though we describe the problem as that of choosing both  $C$  and  $S$  to maximize the objective subject to (A)-(D), effectively we choose only  $C$ , since at each  $t$ , once  $C_t$  is chosen,  $S_{t+1}$  is determined by (2). Note also that (D) means that the formal mathematical problem here is not choosing a sequence of numbers  $\{C_t\}$ , but choosing a sequence of functions

$\{C_t^*(\cdot)\}$  such that at each date  $t$ ,  $C_t^*(w_t) = C_t^*(\{C_{s-1}, S_s, e_s\}, \text{all } s \leq t)$  maps our position  $w_t$  in the information set  $\mathcal{I}_t$  into our best choice for  $C_t$ .<sup>1</sup>

To complete the specification we need assumptions on the random disturbances  $e$ . The standard dynamic programming framework requires that

E) for each  $t$ ,  $e_{t+1}$  is independent of  $C_t$  and of all the random variables in  $\mathcal{I}_t$ <sup>2</sup> and that

F) the random variables  $\{e_t, t = 1, \dots, \infty\}$  are mutually independent and identically distributed (i.i.d.).

This means that no choice made before time  $t$  can influence the realization of the random variable  $e_t$ . Note that this does not mean that  $e_t$  and  $C_s$  are independent for  $s \geq t$ . Since  $e$ 's dated earlier than  $s$  are in  $\mathcal{I}_s$ , they may influence our choice of  $C_s$ ; and this will create dependence in the joint probability distribution of  $C_s$  and  $e_t$  for  $s \geq t$ .

Now observe that the range of probability distributions we can generate for values of  $C_t$  and  $S_{t+1}$  for  $t \geq 0$  through our choice of  $C^*$  functions depends on  $\mathcal{I}_0$  only via  $S_0$ . Since we are allowed to make the choice of  $C_t$  depend on anything in  $\mathcal{I}_t$  that we like, we can create dependencies between the actual future values of  $S$  and  $C$  and, say,  $e_{-3}$  if we like. But since this dependence can take any form we like, and since the data in  $\mathcal{I}_0$  are all fixed and known to us at the time 0 when we choose  $C_0$ , the range of distributions for future  $C$  and  $S$  that we can achieve does not depend on  $\mathcal{I}_0$ , except through the fact that  $S_0$  enters the version of (2) for  $t = 1$ . We denote by  $\mathcal{S}$  the set of all possible values of  $S$ , and we mean by calling  $\mathcal{S}$  “all possible values of  $S$ ” that not only is every value of  $S_0$  with which we might be confronted in  $\mathcal{S}$ , but also for every  $S_0$  in  $\mathcal{S}$  and every way of choosing  $C$ 's that satisfies (A)-(D),  $S_t$  lies in  $\mathcal{S}$  for all  $t$  with probability one. Thus for every  $S_0$  in  $\mathcal{S}$  there will be a unique, possibly infinite, least upper bound for the attainable values of the objective function. [Note that, though the future  $C$ 's and  $S$ 's are unknown and random at time 0, the objective function includes an expectation operator, so its value is a number, not a random variable.] We denote by  $V(\cdot)$  the function mapping  $S$ 's in  $\mathcal{S}$  into the least upper bound of achievable values of the objective function. If the problem is well defined,  $V(S)$  exists for each  $S$  in  $\mathcal{S}$ , though it is important in practice to check that the infinite sum in (1) indeed converges for all feasible choices of actions.  $V$  is called the **value function**.

---

<sup>1</sup> The technically sophisticated reader may note that (D) directly rules out the kind of non-measurability that concerns Stokey and Lucas in their chapter 9.1.

<sup>2</sup> Since  $C_t$  is required by (D) above to depend only on random variables in  $\mathcal{I}_t$ , (E) could omit the “of  $C_t$  and” phrase.

## II. The principle of optimality: necessity and sufficiency

**Theorem 1: (The Principle of Optimality)** Suppose that  $V$  is the value function for the problem of maximizing (1) subject to (A)-(F). Then for each  $S$  in  $\mathcal{S}$ ,  $E[V(f(C, S, e))]$  exists for all  $C$  in  $\Gamma(S)$  (with  $\pm$ infinity allowable values), and

$$V(S) = \underset{C \in \Gamma(S)}{\text{l.u.b.}} \left\{ U(C, S) + bE[V(f(C, S, e))] \right\} \quad (3)$$

*Remark:* In (3) we have omitted dates on  $C$ ,  $S$ , and  $e$ , but we mean here to take the expectation with respect to the distribution of  $e$ , which is the same for all  $t$ , and to treat  $S$  as non-random. If (3) is true in this form for every  $S$  in  $\mathcal{S}$ , then of course it will also be true at every  $t$  with  $C$ ,  $S$  and the  $E$  operator given  $t$  subscripts and  $e$  given a  $t+1$  subscript. For (3) to make sense,  $E[V(C, S, e)]$  must be defined, though possibly infinite, as the theorem asserts.

*Proof:* Note that the objective (1) can be written as

$$U(C_0, S_0) + bE_0 \left[ E_1 \sum_{t=0}^{\infty} b^t U(C_{t+1}, S_{t+1}) \right]. \quad (4)$$

The term in (4) in brackets, together with the preceding  $E$  operator, is exactly the same in form as (1), except with all the time subscripts advanced by 1. Since the constraints are all of the same form at all dates, the least upper bound of this term for a given value of  $S_1$  is  $V(S_1)$ . In a well-defined problem, (4), being the value of the objective function, must itself be well-defined for every  $S_0$  in  $\mathcal{S}$  and every feasible choice of actions. But one particular feasible choice of actions is to choose  $C_0$  in  $\Gamma(S_0)$  arbitrarily, then to choose  $C_t$  for dates  $t=1$  and later so that the second additive term in (4), for every possible value of  $S_1 = f(C_0, S_0, e_1)$ , is at least  $V(S_1) - d$  when  $V(S_1)$  is finite and at least  $1/d$  when  $V(S_1)$  is infinite, where  $d$  is an arbitrarily small positive number. With this particular way of choosing  $C$ 's, we will have, therefore, when  $V(S_1)$  is finite with probability one,

$$\left\{ U(C_0, S_0) + bV(S_1) \right\} - \left\{ U(C_0, S_0) + bE_1 \left[ \sum_{t=0}^{\infty} b^t U(C_{t+1}, S_{t+1}) \right] \right\} \in [0, d]. \quad (5)$$

Since the second term in brackets is a random variable whose expectation we know exists, as it is the value of the objective function for a feasible choice of actions, (5) bounds its first term in brackets above and below by random variables whose expectations exist and are arbitrarily close to each other. Thus the expectation of the first term exists as well. When  $V(S_1)$  is infinite with non-zero probability, our choice of  $C$ 's gives an arbitrarily large value of the objective function, implying that the first term, being bounded below by random variables with arbitrarily large expectation, itself has a well-defined infinite expectation.

Now suppose (3) were not true. Then either there would be, for some  $S$  in  $\mathcal{S}$ , a choice  $\tilde{C}(S)$  of  $C$  making the right-hand-side of (3) exceed  $V(S)$ , or there would be some  $S$  in  $\mathcal{S}$  such that the right-hand side of (3) is bounded away from  $V(S)$  from below. Suppose that (3) fails through

the right-hand side being larger than  $V(\tilde{S})$ , where  $\tilde{S}$  is the particular value of  $S$  at which (3) fails. Now consider this way of choosing  $C$ 's when  $S_0 = \tilde{S}$ : at time 0, choose  $C_0 = \tilde{C}(\tilde{S})$ ; at all later dates, choose  $C$ 's according to a scheme that makes the term in brackets in (4) very close to  $V(S_1)$ . By doing so, we can make (4) as close as we like to the right-hand side of (3). But then we will have succeeded in choosing  $C$ 's in such a way that the objective function value exceeds  $V(S_0)$ , a contradiction.

If (3) fails the other way, so that for some  $\tilde{S}$  the right-hand side of (3) is bounded away from  $V(\tilde{S})$  from below, a parallel argument, again using (4) shows that there is no way to choose  $C$ 's to bring the objective function value arbitrarily close to  $V(\tilde{S})$  when  $S_0 = \tilde{S}$ , again a contradiction with the definition of  $V$ , which completes the proof.

There is an additional necessary condition on the value function that characterizes its long run rate of growth.

*Theorem 2:* Suppose that  $V$  is the value function for the problem of maximizing (1) subject to (A)-(F). Then for every  $d > 0$ , it is possible to choose a policy function  $C_d^*(\cdot)$  such that, for each  $S$  in  $\mathcal{S}$  with  $V(S)$  finite, the value of the objective function attained using  $C_d^*(\cdot)$  is at least  $V(S) - d$ , and the sequence of  $S_t, t = 1, \dots, \infty$  generated by setting  $S_0 = S$  and

$$S_t = f(C_d^*(S_{t-1}), S_{t-1}, e_t), t = 1, \dots, \infty \quad (6)$$

satisfies

$$b^t E_0 V(S_t) \xrightarrow{t \rightarrow \infty} 0. \quad (7)$$

*Remark:* In the (usual) special case where there is a policy function  $C^*(\cdot)$  that actually generates an objective function value equal to  $V(S_0)$  (rather than just arbitrarily close to it), (7) must hold for the  $S$  sequence generated by  $C^*(\cdot)$  from every initial  $S$  for which  $V(S)$  is finite.

*Proof:* First observe that we can certainly find  $C_d^*(\cdot)$ . Since  $V$  satisfies (3) by Theorem 1, we can for each  $S$  that delivers finite  $V(S)$  choose  $C_d^*(S)$  to satisfy

$$U(C_d^*(S), S) + bE[V(f(C_d^*(S), S, e))] \geq V(S) - (1 - b)d. \quad (8)$$

If (8) holds for each  $S$  with finite  $V(S)$ , then we can apply (8) to the term in brackets on its own left-hand side to obtain

$$U(C_d^*(S_0), S_0) + bE_0[U(C_d^*(S_1), S_1) + bV(S_2)] \geq V(S_0) - (1 - b)d - (1 - b)db, \quad (9)$$

where we are assuming that the  $S_t$  sequence is being generated from (6). Repeatedly applying (8) this way will give us the desired conclusion, that the realized value of the objective function using  $C_d^*(\cdot)$  is at least  $V(S_0) - d$ .

This result suggests a method for solving these problems: keep guessing forms for the  $V$  function until we find one that satisfies (3) for all  $S$  in  $\mathcal{S}$ . Since in checking whether (3) is satisfied for all  $S$ , we will ordinarily be finding, for every  $S$ , the  $C^*(S)$  that maximizes the right-hand side, we will have the policy rule  $C^*(\cdot)$  immediately at hand when we have found the right  $V$ .

The problems with this strategy are, first, that it is hopelessly inefficient until we find some systematic way to locate  $V$ 's that might satisfy (3), and, second, that so far we know only that the  $V$  that represents the maximum attainable objective function value -- the value function of the problem -- satisfies (3). We have not yet shown that there cannot be other functions  $V$  that also satisfy (3). It turns out that generally there are other  $V$ 's, besides the actual value function, that satisfy (3), but that we can pick out the true value function by applying some additional side conditions.

*Theorem 3:* Suppose there is a function  $V^*$  that satisfies (3) for every  $S$  in  $\mathcal{S}$  and that in addition

- i) for every  $S$  in  $\mathcal{S}$  there is a value  $C^*(S)$  for  $C$  that attains the maximum on the right-hand side of (3) with  $V = V^*$ ;
- ii) for every  $S_0$  in  $\mathcal{S}$ , if  $S_t$ ,  $t = 1, \dots, \infty$ , is generated from

$$S_{t+1} = f(C^*(t), S_t, e_{t+1}) \quad (10)$$

then

$$b^t E_0[V^*(S_t)] \xrightarrow{t \rightarrow \infty} 0 ; \text{ and} \quad (11)$$

- iii) for any  $V^{**} \neq V^*$  that solves (3) for every  $S$  in  $\mathcal{S}$ , there is some  $d > 0$  such that for any associated  $C_d^{**}$  satisfying

$$U(C_d^{**}(S), S) + bE[V^{**}(f(C_d^{**}(S), S, e))] \geq V^{**}(S) - d \quad (12)$$

for every  $S_0$  in  $\mathcal{S}$ , if  $S_t$ ,  $t = 1, \dots, \infty$ , are generated from

$$S_t = f(C_d^{**}(S_{t-1}), S_{t-1}, e_t), t = 1, \dots, \infty, \text{ then} \quad (13)$$

$$\lim_{t \rightarrow \infty} b^t E_0 V^*(S_t) \geq 0. \quad (14)$$

Then  $V^*$  is the value function for the problem.

*Remark:* The theorem asserts that if the necessary conditions of Theorems 1 and 2 are met for  $V^*$ , then  $V^*$  is the value function unless there is some other solution to (3) that violates (14). It is straightforward, but somewhat tedious, to extend this theorem to the case where no  $C^*$  is available and we instead have to settle for a sequence of  $C_d^*$ 's that approach the upper bound, as we used in Theorem 2.

*Proof:* Suppose there is some solution  $V^{**} \neq V^*$  to (3). If for all  $S$  in  $\mathcal{S}$ ,  $V^{**}(S) \leq V^*(S)$ , with strict inequality for some  $S$ , then  $V^{**}$  can't be the value function. This follows because, since it satisfies (ii),  $V^*$  does represent an attainable value of the objective function for each possible value of its argument, so a policy that delivers a lower  $V^{**}$  instead cannot be optimal. But then suppose that, for some  $\bar{S}$ ,  $V^{**}(\bar{S}) - V^*(\bar{S}) = g > 0$ . Consider the following policy: for  $t = 0, \dots, T-1$ , set  $C_t = C_d^{**}(S_t)$  chosen as in (iii), then for  $t \geq T$  set  $C_t = C^*(S_t)$ , . The value of the objective function under this policy starting from  $S_0 = \bar{S}$  is

$$E_0 \left[ \sum_{t=0}^{T-1} U(C_d^{**}(S_t), S_t) b^t \right] + b^T E_0 V^*(S_T) \xrightarrow{T \rightarrow \infty} V_d^{**}(\bar{S}), \quad (15)$$

where by making  $d$  small enough we can make  $V_d^{**}(\bar{S})$  as close as we like to  $V^{**}(\bar{S})$ . The convergence in (15) follows from (iii) and the definition of  $C_d^{**}$ . But notice that

$$V^*(\bar{S}) = U(C^*(\bar{S}), \bar{S}) + b E_0 V^*(S_1) \geq U(C_d^{**}(\bar{S}), \bar{S}) + b E_0 V^*(S_1) \quad (16)$$

(Note that in (16) we are implicitly treating the two  $S_1$  values as generated by the  $C$  choice at time zero in each expression -- so the  $S_1$ 's in the two expressions are not the same.) Applying the same argument again to  $V^{**}(S_1)$  in (16), and so on recursively  $T$  times allows us to conclude

$$E_0 \left[ \sum_{t=0}^{T-1} U(C_d^{**}(S_t), S_t) b^t \right] + b^T E_0 V^*(S_T) < V^*(\bar{S}). \quad (17)$$

But with  $d$  arbitrarily small, (15) and (17) together imply  $V^*(\bar{S}) > V^{**}(\bar{S})$ , which contradicts our initial assumption. So  $V^{**}(S) \leq V^*(S)$  for all  $S$ , which we have already noted means that  $V^{**}$  is not the value function. Since the argument applies to arbitrary  $V^{**} \neq V^*$ ,  $V^*$  is the unique optimal solution and we have completed the proof.

*Corollary:* If  $U$  is bounded below, the necessary conditions of Theorems 1 and 2 are also sufficient.

*Proof:* Because the objective function is discounted, it is bounded below when  $U$  is bounded below. This makes (14) hold automatically.

### III. Value Iteration

The arguments of the previous section point to a conceptually simple method for approximating the value function and hence  $C^*$ , the optimal policy rule. The method is called, for reasons that will be obvious, **value function iteration**, and proceeds as follows. Begin by guessing a form  $V_0$  for the value function. Then for each iteration  $n$ ,  $n=1,2,3,\dots$ , set, for each  $S$  in  $\mathcal{S}$ ,

$$V_n(S) = \text{l.u.b.}_{C \text{ in } \Gamma} \{ U(C, S) + b E[V_{n-1}(f(C, S, e))] \}. \quad (18)$$

Continue until  $V_n(S) = V_{n-1}(S)$ , all  $S$ , to within some criterion for numerical accuracy. At that point a  $V$  satisfying the principal of optimality will have been arrived at, and the other necessary and sufficient conditions can be checked.

It is not necessarily true that value function iterations converge, however. When  $U$  is unbounded either above or below, it can easily happen that value function iteration convergence fails. Since many of the standard utility functions of macroeconomic models -- logarithmic

$U(C) = \log(C)$  and CRR  $U(C) = \frac{C^{1-g}}{1-g}$  for example -- fail to satisfy such a boundedness

condition, we must generally be wary that value iteration might not converge. It is worth knowing, though, that when  $U$  is bounded (which we already know is a sufficient condition to guarantee that a solution to the optimality equation is the value function) value iteration necessarily converges. The argument goes as follows. First we note that

$$|V_n(S) - V_{n-1}(S)| \leq \left| \text{l.u.b.}_{C \in \Gamma(S)} \left\{ U(C, S) - U(C, S) + bE[V_{n-1}(f(C, S, e)) - V_{n-2}(f(C, S, e))] \right\} \right|. \quad (19)$$

This follows from (18) and the fact that

$$|\text{l.u.b.}\{a(\cdot)\} - \text{l.u.b.}\{b(\cdot)\}| \leq \text{l.u.b.}\{|a(\cdot) - b(\cdot)|\}. \quad (20)$$

If we introduce as a norm on the space of  $V$ 's

$$\|V\| = \sup_S |V(S)|, \quad (21)$$

(19) can be used to produce

$$\|V_n - V_{n-1}\| \leq b\|V_{n-1} - V_{n-2}\|. \quad (22)$$

This implies that the sequence of value function iterates  $\{V_n\}$  is a Cauchy sequence on the space of bounded functions, and hence that it converges to some bounded function. The reason the argument fails when  $U$  is unbounded is that then  $V$  is generally unbounded (not infinite -- it just gets arbitrarily large or small as we change its argument  $S$ ) and therefore does not allow us to use (21). Even if we start with a bounded  $V_0$ , the unboundedness of  $U$  generally can make  $V_1$  unbounded.