

Dynamic Programming and Optimal Control

Richard Weber

Graduate Course at London Business School

Autumn 2014

Contents

Table of Contents	ii
1 Dynamic Programming	1
1.1 Control as optimization over time	1
1.2 The principle of optimality	1
1.3 Example: the shortest path problem	1
1.4 The optimality equation	2
1.5 Markov decision processes	4
2 Examples of Dynamic Programming	5
2.1 Example: optimization of consumption	5
2.2 Example: exercising a stock option	6
2.3 Example: secretary problem	7
3 Dynamic Programming over the Infinite Horizon	9
3.1 Discounted costs	9
3.2 Example: job scheduling	9
3.3 The infinite-horizon case	10
3.4 The optimality equation in the infinite-horizon case	11
3.5 Example: selling an asset	12
4 Positive Programming	14
4.1 Example: possible lack of an optimal policy.	14
4.2 Characterization of the optimal policy	14
4.3 Example: optimal gambling	15
4.4 Value iteration	15
4.5 Example: search for a moving object	16
4.6 Example: pharmaceutical trials	17
5 Negative Programming	19
5.1 Example: a partially observed MDP	19
5.2 Stationary policies	20
5.3 Characterization of the optimal policy	20
5.4 Optimal stopping over a finite horizon	21
5.5 Example: optimal parking	22
6 Optimal Stopping Problems	23
6.1 Bruss's odds algorithm	23
6.2 Example: Stopping a random walk	24
6.3 Optimal stopping over the infinite horizon	24
6.4 Sequential Probability Ratio Test	26
6.5 Bandit processes	26
6.6 Example: Two-armed bandit	27

6.7	Example: prospecting	27
7	Bandit Processes and the Gittins Index	29
7.1	Index policies	29
7.2	Multi-armed bandit problem	29
7.3	Gittins index theorem	30
7.4	Calibration	31
7.5	Proof of the Gittins index theorem	31
7.6	Example: Weitzman's problem	32
8	Applications of Bandit Processes	33
8.1	Forward induction policies	33
8.2	Example: playing golf with more than one ball	33
8.3	Target processes	34
8.4	Bandit superprocesses	34
8.5	Example: single machine stochastic scheduling	35
8.6	Calculation of the Gittins index	35
8.7	Branching bandits	36
8.8	Example: Searching for a single object	37
9	Average-cost Programming	38
9.1	Average-cost optimality equation	38
9.2	Example: admission control at a queue	39
9.3	Value iteration bounds	39
9.4	Policy improvement algorithm	40
10	Continuous-time Markov Decision Processes	42
10.1	Stochastic scheduling on parallel machines	42
10.2	Controlled Markov jump processes	44
10.3	Example: admission control at a queue	45
11	Restless Bandits	47
11.1	Examples	47
11.2	Whittle index policy	48
11.3	Whittle indexability	49
11.4	Fluid models of large stochastic systems	49
11.5	Asymptotic optimality	50
12	Sequential Assignment and Allocation Problems	53
12.1	Sequential stochastic assignment problem	53
12.2	Sequential allocation problems	54
12.3	SSAP with arrivals	56
12.4	SSAP with a postponement option	57
12.5	Stochastic knapsack and bin packing problems	58

13 LQ Regulation	59
13.1 The LQ regulation problem	59
13.2 The Riccati recursion	61
13.3 White noise disturbances	61
13.4 LQ regulation in continuous-time	62
13.5 Linearization of nonlinear models	62
14 Controllability and Observability	63
14.1 Controllability and Observability	63
14.2 Controllability	63
14.3 Controllability in continuous-time	65
14.4 Example: broom balancing	65
14.5 Stabilizability	66
14.6 Example: pendulum	66
14.7 Example: satellite in a plane orbit	67
15 Observability and the LQG Model	68
15.1 Infinite horizon limits	68
15.2 Observability	68
15.3 Observability in continuous-time	70
15.4 Example: satellite in planar orbit	70
15.5 Imperfect state observation with noise	70
16 Kalman Filter and Certainty Equivalence	72
16.1 The Kalman filter	72
16.2 Certainty equivalence	73
16.3 The Hamilton-Jacobi-Bellman equation	74
16.4 Example: LQ regulation	75
16.5 Example: harvesting fish	75
17 Pontryagin's Maximum Principle	78
17.1 Example: optimization of consumption	78
17.2 Heuristic derivation of Pontryagin's maximum principle	79
17.3 Example: parking a rocket car	80
17.4 Adjoint variables as Lagrange multipliers	82
18 Using Pontryagin's Maximum Principle	83
18.1 Transversality conditions	83
18.2 Example: use of transversality conditions	83
18.3 Example: insects as optimizers	84
18.4 Problems in which time appears explicitly	84
18.5 Example: monopolist	85
18.6 Example: neoclassical economic growth	86

19 Controlled Diffusion Processes	88
19.1 The dynamic programming equation	88
19.2 Diffusion processes and controlled diffusion processes	88
19.3 Example: noisy LQ regulation in continuous time	89
19.4 Example: passage to a stopping set	90
20 Risk-sensitive Optimal Control	92
20.1 Whittle risk sensitivity	92
20.2 The felicity of LEQG assumptions	92
20.3 A risk-sensitive certainty-equivalence principle	94
20.4 Large deviations	95
20.5 A risk-sensitive maximum principle	96
20.6 Example: risk-sensitive optimization of consumption	96
Index	97

1 Dynamic Programming

Dynamic programming and the principle of optimality. Notation for state-structured models. Feedback, open-loop, and closed-loop controls. Markov decision processes.

1.1 Control as optimization over time

Optimization is a key tool in modelling. Sometimes it is important to solve a problem optimally. Other times either a near-optimal solution is good enough, or the real problem does not have a single criterion by which a solution can be judged. However, even when an optimal solution is not required it can be useful to test one's thinking by following an optimization approach. If the 'optimal' solution is ridiculous it may suggest ways in which both modelling and thinking can be refined.

Control theory is concerned with dynamic systems and their **optimization over time**. It accounts for the fact that a dynamic system may evolve stochastically and that key variables may be unknown or imperfectly observed.

The optimization models in the IB course (for linear programming and network flow models) were static and nothing was either random or hidden. In this course it is the additional features of dynamic and stochastic evolution, and imperfect state observation, that give rise to new types of optimization problem and which require new ways of thinking.

We could spend an entire lecture discussing the importance of control theory and tracing its development through the windmill, steam governor, and so on. Such 'classic control theory' is largely concerned with the question of stability, and there is much of this theory which we ignore, e.g., Nyquist criterion and dynamic lags.

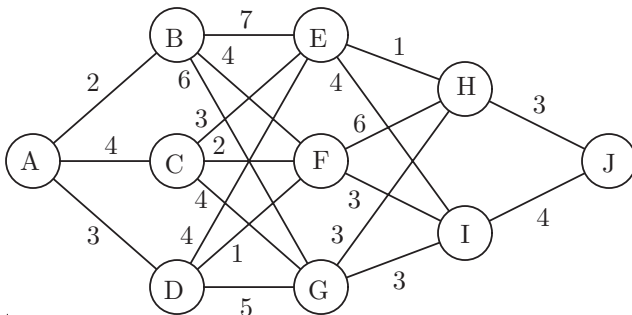
1.2 The principle of optimality

A key idea in this course is that optimization over time can often be seen as 'optimization in stages'. We trade off our desire to obtain the least possible cost at the present stage against the implication this would have for costs at future stages. The best action minimizes the sum of the cost incurred at the current stage and the least total cost that can be incurred from all subsequent stages, consequent on this decision. This is known as the Principle of Optimality.

Definition 1.1 (Principle of Optimality). From any point on an optimal trajectory, the remaining trajectory is optimal for the problem initiated at that point.

1.3 Example: the shortest path problem

Consider the 'stagecoach problem' in which a traveler wishes to minimize the length of a journey from town A to town J by first traveling to one of B, C or D and then onwards to one of E, F or G then onwards to one of H or I and the finally to J. Thus there are 4 'stages'. The arcs are marked with distances between towns.



Road system for stagecoach problem

Solution. Let $F(X)$ be the minimal distance required to reach J from X. Then clearly, $F(J) = 0$, $F(H) = 3$ and $F(I) = 4$.

$$F(F) = \min[6 + F(H), 3 + F(I)] = 7,$$

and so on. Recursively, we obtain $F(A) = 11$ and simultaneously an optimal route, i.e. $A \rightarrow D \rightarrow F \rightarrow I \rightarrow J$ (although it is not unique).

The study of dynamic programming dates from Richard Bellman, who wrote the first book on the subject (1957) and gave it its name. A very large number of problems can be treated this way.

1.4 The optimality equation

The optimality equation in the general case. In **discrete-time** t takes integer values, say $t = 0, 1, \dots$. Suppose u_t is a **control variable** whose value is to be chosen at time t . Let $U_{t-1} = (u_0, \dots, u_{t-1})$ denote the partial sequence of controls (or decisions) taken over the first t stages. Suppose the cost up to the **time horizon** h is given by

$$\mathbf{C} = G(U_{h-1}) = G(u_0, u_1, \dots, u_{h-1}).$$

Then the **principle of optimality** is expressed in the following theorem.

Theorem 1.2 (The principle of optimality). *Define the functions*

$$G(U_{t-1}, t) = \inf_{u_t, u_{t+1}, \dots, u_{h-1}} G(U_{h-1}).$$

Then these obey the recursion

$$G(U_{t-1}, t) = \inf_{u_t} G(U_t, t+1) \quad t < h,$$

with terminal evaluation $G(U_{h-1}, h) = G(U_{h-1})$.

The proof is immediate from the definition of $G(U_{t-1}, t)$, i.e.

$$G(U_{t-1}, t) = \inf_{u_t} \left\{ \inf_{u_{t+1}, \dots, u_{h-1}} G(u_0, \dots, u_{t-1}, u_t, u_{t+1}, \dots, u_{h-1}) \right\}.$$

The state structured case. The control variable u_t is chosen on the basis of knowing $U_{t-1} = (u_0, \dots, u_{t-1})$, (which determines everything else). But a more economical representation of the past history is often sufficient. For example, we may not need to know the entire path that has been followed up to time t , but only the place to which it has taken us. The idea of a **state variable** $x \in \mathbb{R}^d$ is that its value at t , denoted x_t , can be found from known quantities and obeys a **plant equation** (or law of motion)

$$x_{t+1} = a(x_t, u_t, t).$$

Suppose we wish to minimize a **separable cost function** of the form

$$\mathbf{C} = \sum_{t=0}^{h-1} c(x_t, u_t, t) + \mathbf{C}_h(x_h), \quad (1.1)$$

by choice of controls $\{u_0, \dots, u_{h-1}\}$. Define the cost from time t onwards as,

$$\mathbf{C}_t = \sum_{\tau=t}^{h-1} c(x_\tau, u_\tau, \tau) + \mathbf{C}_h(x_h), \quad (1.2)$$

and the minimal cost from time t onwards as an optimization over $\{u_t, \dots, u_{h-1}\}$ conditional on $x_t = x$,

$$F(x, t) = \inf_{u_t, \dots, u_{h-1}} \mathbf{C}_t.$$

Here $F(x, t)$ is the minimal future cost from time t onward, given that the state is x at time t . Then by an inductive proof, one can show as in Theorem 1.2 that

$$F(x, t) = \inf_u [c(x, u, t) + F(a(x, u, t), t + 1)], \quad t < h, \quad (1.3)$$

with terminal condition $F(x, h) = \mathbf{C}_h(x)$. Here x is a generic value of x_t . The minimizing u in (1.3) is the optimal control $u(x, t)$ and values of x_0, \dots, x_{t-1} are irrelevant. The **optimality equation** (1.3) is also called the **dynamic programming equation** (DP) or **Bellman equation**.

The DP equation defines an optimal control problem in what is called **feedback** or **closed-loop** form, with $u_t = u(x_t, t)$. This is in contrast to the **open-loop** formulation in which $\{u_0, \dots, u_{h-1}\}$ are to be determined all at once at time 0. A **policy** (or strategy) is a rule for choosing the value of the control variable under all possible circumstances as a function of the perceived circumstances. To summarise:

- (i) The optimal u_t is a function only of x_t and t , i.e. $u_t = u(x_t, t)$.
- (ii) The DP equation expresses the optimal u_t in closed-loop form. It is optimal whatever the past control policy may have been.
- (iii) The DP equation is a backward recursion in time (from which we get the optimum at $h - 1$, then $h - 2$ and so on.) The later policy is decided first.

‘Life must be lived forward and understood backwards.’ (Kierkegaard)

1.5 Markov decision processes

Consider now stochastic evolution. Let $X_t = (x_0, \dots, x_t)$ and $U_t = (u_0, \dots, u_t)$ denote the x and u histories at time t . As above, state structure is characterized by the fact that the evolution of the process is described by a state variable x , having value x_t at time t . The following assumptions define what is known as a discrete-time **Markov decision process** (MDP).

(a) *Markov dynamics*: (i.e. the stochastic version of the plant equation.)

$$P(x_{t+1} \mid X_t, U_t) = P(x_{t+1} \mid x_t, u_t).$$

(b) *Separable (or decomposable) cost function*, (i.e. cost given by (1.1)).

For the moment we also require the following:

(c) *Perfect state observation*: The current value of the state is observable. That is, x_t is known when choosing u_t . So, letting W_t denote the observed history at time t , we assume $W_t = (X_t, U_{t-1})$.

Note that \mathbf{C} is determined by W_h , so we might write $\mathbf{C} = \mathbf{C}(W_h)$.

As in the previous section, the cost from time t onwards is given by (1.2). Denote the minimal expected cost from time t onwards by

$$F(W_t) = \inf_{\pi} E_{\pi}[\mathbf{C}_t \mid W_t],$$

where π denotes a policy, i.e. a rule for choosing the controls u_0, \dots, u_{h-1} .

The following theorem is then obvious.

Theorem 1.3. *$F(W_t)$ is a function of x_t and t alone, say $F(x_t, t)$. It obeys the optimality equation*

$$F(x_t, t) = \inf_{u_t} \{c(x_t, u_t, t) + E[F(x_{t+1}, t+1) \mid x_t, u_t]\}, \quad t < h, \quad (1.4)$$

with terminal condition

$$F(x_h, h) = \mathbf{C}_h(x_h).$$

Moreover, a minimizing value of u_t in (1.4) (which is also only a function x_t and t) is optimal.

Proof. The value of $F(W_h)$ is $\mathbf{C}_h(x_h)$, so the asserted reduction of F is valid at time h . Assume it is valid at time $t+1$. The DP equation is then

$$F(W_t) = \inf_{u_t} \{c(x_t, u_t, t) + E[F(x_{t+1}, t+1) \mid X_t, U_t]\}. \quad (1.5)$$

But, by assumption (a), the right-hand side of (1.5) reduces to the right-hand member of (1.4). All the assertions then follow. \square

2 Examples of Dynamic Programming

Examples of dynamic programming problems and some useful tricks to solve them. The idea that it can be useful to model things in terms of time to go.

2.1 Example: optimization of consumption

An investor receives annual income of x_t pounds in year t . He consumes u_t and adds $x_t - u_t$ to his capital, $0 \leq u_t \leq x_t$. The capital is invested at interest rate $\theta \times 100\%$, and so his income in year $t + 1$ increases to

$$x_{t+1} = a(x_t, u_t) = x_t + \theta(x_t - u_t). \quad (2.1)$$

He desires to maximize total consumption over h years,

$$\mathbf{C} = \sum_{t=0}^{h-1} c(x_t, u_t, t) + \mathbf{C}_h(x_h) = \sum_{t=0}^{h-1} u_t$$

The plant equation (2.1) specifies a **Markov decision process** (MDP). When we add to this the aim of maximizing the performance measure \mathbf{C} we have what is called a **Markov decision problem**. For both we use the abbreviation MDP. In the notation we have been using, $c(x_t, u_t, t) = u_t$, $\mathbf{C}_h(x_h) = 0$. This is termed a **time-homogeneous** model because neither costs nor dynamics depend on t .

Solution. Since dynamic programming makes its calculations backwards, from the termination point, it is often advantageous to write things in terms of the ‘**time to go**’, $s = h - t$. Let $F_s(x)$ denote the maximal reward obtainable, starting in state x when there is time s to go. The dynamic programming equation is

$$F_s(x) = \max_{0 \leq u \leq x} [u + F_{s-1}(x + \theta(x - u))],$$

where $F_0(x) = 0$, (since nothing more can be consumed once time h is reached.) Here, x and u are generic values for x_s and u_s .

We can substitute backwards and soon guess the form of the solution. First,

$$F_1(x) = \max_{0 \leq u \leq x} [u + F_0(u + \theta(x - u))] = \max_{0 \leq u \leq x} [u + 0] = x.$$

Next,

$$F_2(x) = \max_{0 \leq u \leq x} [u + F_1(x + \theta(x - u))] = \max_{0 \leq u \leq x} [u + x + \theta(x - u)].$$

Since $u + x + \theta(x - u)$ linear in u , its maximum occurs at $u = 0$ or $u = x$, and so

$$F_2(x) = \max[(1 + \theta)x, 2x] = \max[1 + \theta, 2]x = \rho_2 x.$$

This motivates the guess $F_{s-1}(x) = \rho_{s-1}x$. Trying this, we find

$$F_s(x) = \max_{0 \leq u \leq x} [u + \rho_{s-1}(x + \theta(x - u))] = \max[(1 + \theta)\rho_{s-1}, 1 + \rho_{s-1}]x = \rho_s x.$$

Thus our guess is verified and $F_s(x) = \rho_s x$, where ρ_s obeys the recursion implicit in the above, and i.e. $\rho_s = \rho_{s-1} + \max[\theta\rho_{s-1}, 1]$. This gives

$$\rho_s = \begin{cases} s & s \leq s^* \\ (1 + \theta)^{s-s^*} s^* & s \geq s^* \end{cases},$$

where s^* is the least integer such that $1 + s^* \leq (1 + \theta)s^* \iff s^* \geq 1/\theta$, i.e. $s^* = \lceil 1/\theta \rceil$. The optimal strategy is to invest the whole of the income in years $0, \dots, h - s^* - 1$, (to build up capital) and then consume the whole of the income in years $h - s^*, \dots, h - 1$.

There are several things worth learning from this example. (i) It is often useful to frame things in terms of time to go, s . (ii) Although the form of the dynamic programming equation can sometimes look messy, try working backwards from $F_0(x)$ (which is known). Often a pattern will emerge from which you can piece together a solution. (iii) When the dynamics are linear, the optimal control lies at an extreme point of the set of feasible controls. This form of policy, which either consumes nothing or consumes everything, is known as **bang-bang control**.

2.2 Example: exercising a stock option

The owner of a call option has the option to buy a share at fixed ‘striking price’ p . The option must be exercised by day h . If she exercises the option on day t and then immediately sells the share at the current price x_t , she can make a profit of $x_t - p$. Suppose the price sequence obeys the equation $x_{t+1} = x_t + \epsilon_t$, where the ϵ_t are i.i.d. random variables for which $E|\epsilon| < \infty$. The aim is to exercise the option optimally.

Let $F_s(x)$ be the **value function** (maximal expected profit) when the share price is x and there are s days to go. Show that (i) $F_s(x)$ is non-decreasing in s , (ii) $F_s(x) - x$ is non-increasing in x and (iii) $F_s(x)$ is continuous in x . Deduce that the optimal policy can be characterized as follows.

There exists a non-decreasing sequence $\{a_s\}$ such that an optimal policy is to exercise the option the first time that $x \geq a_s$, where x is the current price and s is the number of days to go before expiry of the option.

Solution. The state variable at time t is, strictly speaking, x_t plus a variable which indicates whether the option has been exercised or not. However, it is only the latter case which is of interest, so x is the effective state variable. As above, we use time to go, $s = h - t$. So if we let $F_s(x)$ be the value function (maximal expected profit) with s days to go then

$$F_0(x) = \max\{x - p, 0\},$$

and so the dynamic programming equation is

$$F_s(x) = \max\{x - p, E[F_{s-1}(x + \epsilon)]\}, \quad s = 1, 2, \dots$$

Note that the expectation operator comes *outside*, not inside, $F_{s-1}(\cdot)$.

It is easy to show (i), (ii) and (iii) by induction on s . For example, (i) is obvious, since increasing s means we have more time over which to exercise the option. However, for a formal proof

$$F_1(x) = \max\{x - p, E[F_0(x + \epsilon)]\} \geq \max\{x - p, 0\} = F_0(x).$$

Now suppose, inductively, that $F_{s-1} \geq F_{s-2}$. Then

$$F_s(x) = \max\{x - p, E[F_{s-1}(x + \epsilon)]\} \geq \max\{x - p, E[F_{s-2}(x + \epsilon)]\} = F_{s-1}(x),$$

whence F_s is non-decreasing in s . Similarly, an inductive proof of (ii) follows from

$$\underbrace{F_s(x) - x}_{\geq 0} = \max\{-p, \underbrace{E[F_{s-1}(x + \epsilon) - (x + \epsilon)]}_{\leq 0} + E(\epsilon)\},$$

since the left hand underbraced term inherits the non-increasing character of the right hand underbraced term. Thus the optimal policy can be characterized as stated. For from (ii), (iii) and the fact that $F_s(x) \geq x - p$ it follows that there exists an a_s such that $F_s(x)$ is greater than $x - p$ if $x < a_s$ and equals $x - p$ if $x \geq a_s$. It follows from (i) that a_s is non-decreasing in s . The constant a_s is the smallest x for which $F_s(x) = x - p$.

2.3 Example: secretary problem

We are to interview h candidates for a job. At the end of each interview we must either hire or reject the candidate we have just seen, and may not change this decision later. Candidates are seen in random order and can be ranked against those seen previously. The aim is to maximize the probability of choosing the candidate of greatest rank.

Solution. Let W_t be the history of observations up to time t , i.e. after we have interviewed the t th candidate. All that matters are the value of t and whether the t th candidate is better than all her predecessors: let $x_t = 1$ if this is true and $x_t = 0$ if it is not. In the case $x_t = 1$, the probability she is the best of all h candidates is

$$P(\text{best of } h \mid \text{best of first } t) = \frac{P(\text{best of } h)}{P(\text{best of first } t)} = \frac{1/h}{1/t} = \frac{t}{h}.$$

Now the fact that the t th candidate is the best of the t candidates seen so far places no restriction on the relative ranks of the first $t - 1$ candidates; thus $x_t = 1$ and W_{t-1} are statistically independent and we have

$$P(x_t = 1 \mid W_{t-1}) = \frac{P(W_{t-1} \mid x_t = 1)}{P(W_{t-1})} P(x_t = 1) = P(x_t = 1) = \frac{1}{t}.$$

Let $F(t - 1)$ be the probability that under an optimal policy we select the best candidate, given that we have passed over the first $t - 1$ candidates. Dynamic programming gives

$$F(t-1) = \frac{t-1}{t}F(t) + \frac{1}{t} \max\left(\frac{t}{h}, F(t)\right) = \max\left(\frac{t-1}{t}F(t) + \frac{1}{h}, F(t)\right)$$

The first term deals with what happens when the t th candidate is not the best so far; we should certainly pass over her. The second term deals with what happens when she is the best so far. Now we have a choice: either accept her (and she will turn out to be best with probability t/h), or pass over her.

These imply $F(t-1) \geq F(t)$ for all $t \leq h$. Therefore, since t/h and $F(t)$ are respectively increasing and non-increasing in t , it must be that for small t we have $F(t) > t/h$ and for large t we have $F(t) \leq t/h$. Let t_0 be the smallest t such that $F(t) \leq t/h$. Then

$$F(t-1) = \begin{cases} F(t_0), & t < t_0, \\ \frac{t-1}{t}F(t) + \frac{1}{h}, & t \geq t_0. \end{cases}$$

Solving the second of these backwards from the point $t = h$, $F(h) = 0$, we obtain

$$\frac{F(t-1)}{t-1} = \frac{1}{h(t-1)} + \frac{F(t)}{t} = \dots = \frac{1}{h(t-1)} + \frac{1}{ht} + \dots + \frac{1}{h(h-1)},$$

whence

$$F(t-1) = \frac{t-1}{h} \sum_{\tau=t-1}^{h-1} \frac{1}{\tau}, \quad t \geq t_0.$$

Since we require $F(t_0) \leq t_0/h$, it must be that t_0 is the smallest integer satisfying

$$\sum_{\tau=t_0}^{h-1} \frac{1}{\tau} \leq 1.$$

For large h the sum on the left above is about $\log(h/t_0)$, so $\log(h/t_0) \approx 1$ and we find $t_0 \approx h/e$. Thus the optimal policy is to interview $\approx h/e$ candidates, but without selecting any of these, and then select the first candidate thereafter who is the best of all those seen so far. The probability of success is $F(0) = F(t_0) \sim t_0/h \sim 1/e = 0.3679$. It is surprising that the probability of success is so large for arbitrarily large h .

There are a couple things to learn from this example. (i) It is often useful to try to establish the fact that terms over which a maximum is being taken are monotone in opposite directions, as we did with t/h and $F(t)$. (ii) A typical approach is to first determine the form of the solution, then find the optimal cost (reward) function by backward recursion from the terminal point, where its value is known.

3 Dynamic Programming over the Infinite Horizon

Cases of discounted, negative and positive dynamic programming. Validity of the optimality equation over the infinite horizon.

3.1 Discounted costs

For a **discount factor**, $\beta \in (0, 1]$, the **discounted-cost criterion** is defined as

$$\mathbf{C} = \sum_{t=0}^{h-1} \beta^t c(x_t, u_t, t) + \beta^h \mathbf{C}_h(x_h). \quad (3.1)$$

This simplifies things mathematically, particularly when we want to consider an infinite horizon. If costs are uniformly bounded, say $|c(x, u)| < B$, and discounting is strict ($\beta < 1$) then the infinite horizon cost is bounded by $B/(1 - \beta)$. In finance, if there is an interest rate of $r\%$ per unit time, then a unit amount of money at time t is worth $\rho = 1 + r/100$ at time $t + 1$. Equivalently, a unit amount at time $t + 1$ has present value $\beta = 1/\rho$. The function, $F(x, t)$, which expresses the minimal present value at time t of expected-cost from time t up to h is

$$F(x, t) = \inf_{\pi} E_{\pi} \left[\sum_{\tau=t}^{h-1} \beta^{\tau-t} c(x_{\tau}, u_{\tau}, \tau) + \beta^{h-t} \mathbf{C}_h(x_h) \mid x_t = x \right]. \quad (3.2)$$

where E_{π} denotes expectation over the future path of the process under policy π . The DP equation is now

$$F(x, t) = \inf_u [c(x, u, t) + \beta E F(x_{t+1}, t + 1)], \quad t < h, \quad (3.3)$$

where $F(x, h) = \mathbf{C}_h(x)$.

3.2 Example: job scheduling

A collection of n jobs is to be processed in arbitrary order by a single machine. Job i has processing time p_i and when it completes a reward r_i is obtained. Find the order of processing that maximizes the sum of the discounted rewards.

Solution. Here we take ‘time-to-go k ’ as the point at which the $n - k$ th job has just been completed and there remains a set of k uncompleted jobs, say S_k . The dynamic programming equation is

$$F_k(S_k) = \max_{i \in S_k} [r_i \beta^{p_i} + \beta^{p_i} F_{k-1}(S_k - \{i\})].$$

Obviously $F_0(\emptyset) = 0$. Applying the method of dynamic programming we first find $F_1(\{i\}) = r_i \beta^{p_i}$. Then, working backwards, we find

$$F_2(\{i, j\}) = \max[r_i \beta^{p_i} + \beta^{p_i+p_j} r_j, r_j \beta^{p_j} + \beta^{p_j+p_i} r_i].$$

There will be 2^n equations to evaluate, but with perseverance we can determine $F_n(\{1, 2, \dots, n\})$. However, there is a simpler way.

An interchange argument

Suppose jobs are processed in the order $i_1, \dots, i_k, i, j, i_{k+3}, \dots, i_n$. Compare the reward that is obtained if the order of jobs i and j is reversed: $i_1, \dots, i_k, j, i, i_{k+3}, \dots, i_n$. The rewards under the two schedules are respectively

$$R_1 + \beta^{T+p_i} r_i + \beta^{T+p_i+p_j} r_j + R_2 \quad \text{and} \quad R_1 + \beta^{T+p_j} r_j + \beta^{T+p_j+p_i} r_i + R_2,$$

where $T = p_{i_1} + \dots + p_{i_k}$, and R_1 and R_2 are respectively the sum of the rewards due to the jobs coming before and after jobs i, j ; these are the same under both schedules. The reward of the first schedule is greater if $r_i \beta^{p_i} / (1 - \beta^{p_i}) > r_j \beta^{p_j} / (1 - \beta^{p_j})$. Hence a schedule can be optimal only if the jobs are taken in decreasing order of the indices $r_i \beta^{p_i} / (1 - \beta^{p_i})$. This type of reasoning is known as an **interchange argument**.

There are a couple points to note. (i) An interchange argument can be useful for solving a decision problem about a system that evolves in stages. Although such problems can be solved by dynamic programming, an interchange argument – when it works – is usually easier. (ii) The decision points need not be equally spaced in time. Here they are the times at which jobs complete.

3.3 The infinite-horizon case

In the finite-horizon case the value function is obtained simply from (3.3) by the backward recursion from the terminal point. However, when the horizon is infinite there is no terminal point and so the validity of the optimality equation is no longer obvious.

Consider the time-homogeneous Markov case, in which costs and dynamics do not depend on t , i.e. $c(x, u, t) = c(x, u)$. Suppose also that there is no terminal cost, i.e. $C_h(x) = 0$. Define the s -horizon cost under policy π as

$$F_s(\pi, x) = E_\pi \left[\sum_{t=0}^{s-1} \beta^t c(x_t, u_t) \mid x_0 = x \right],$$

If we take the infimum with respect to π we have the *infimal s -horizon cost*

$$F_s(x) = \inf_{\pi} F_s(\pi, x).$$

Clearly, this always exists and satisfies the optimality equation

$$F_s(x) = \inf_u \{c(x, u) + \beta E[F_{s-1}(x_1) \mid x_0 = x, u_0 = u]\}, \quad (3.4)$$

with terminal condition $F_0(x) = 0$.

Sometimes a nice way to write (3.4) is as $F_s = \mathcal{L}F_{s-1}$ where \mathcal{L} is the operator with action

$$\mathcal{L}\phi(x) = \inf_u \{c(x, u) + \beta E[\phi(x_1) \mid x_0 = x, u_0 = u]\}.$$

This operator transforms a scalar function of the state x to another scalar function of x . Note that \mathcal{L} is a **monotone operator**, in the sense that if $\phi_1 \leq \phi_2$ then $\mathcal{L}\phi_1 \leq \mathcal{L}\phi_2$.

The *infinite-horizon cost under policy* π is also quite naturally defined as

$$F(\pi, x) = \lim_{s \rightarrow \infty} F_s(\pi, x). \quad (3.5)$$

This limit need not exist (e.g. if $\beta = 1$, $x_{t+1} = -x_t$ and $c(x, u) = x$), but it will do so under any of the following three scenarios.

- D (**discounted programming**): $0 < \beta < 1$, and $|c(x, u)| < B$ for all x, u .
- N (**negative programming**): $0 < \beta \leq 1$, and $c(x, u) \geq 0$ for all x, u .
- P (**positive programming**): $0 < \beta \leq 1$, and $c(x, u) \leq 0$ for all x, u .

Notice that the names ‘negative’ and ‘positive’ appear to be the wrong way around with respect to the sign of $c(x, u)$. The names actually come from equivalent problems of maximizing rewards, like $r(x, u) (= -c(x, u))$. Maximizing positive rewards (P) is the same thing as minimizing negative costs. Maximizing negative rewards (N) is the same thing as minimizing positive costs. In cases N and P we usually take $\beta = 1$.

The existence of the limit (possibly infinite) in (3.5) is assured in cases N and P by monotone convergence, and in case D because the total cost occurring after the s th step is bounded by $\beta^s B / (1 - \beta)$.

3.4 The optimality equation in the infinite-horizon case

The *infimal infinite-horizon cost* is defined as

$$F(x) = \inf_{\pi} F(\pi, x) = \inf_{\pi} \lim_{s \rightarrow \infty} F_s(\pi, x). \quad (3.6)$$

The following theorem justifies our writing the optimality equation (i.e. (3.7)).

Theorem 3.1. *Suppose D, N, or P holds. Then $F(x)$ satisfies the optimality equation*

$$F(x) = \inf_u \{c(x, u) + \beta E[F(x_1) \mid x_0 = x, u_0 = u]\}. \quad (3.7)$$

Proof. We first prove that ‘ \geq ’ holds in (3.7). Suppose π is a policy, which chooses $u_0 = u$ when $x_0 = x$. Then

$$F_s(\pi, x) = c(x, u) + \beta E[F_{s-1}(\pi, x_1) \mid x_0 = x, u_0 = u]. \quad (3.8)$$

Either D, N or P is sufficient to allow us to take limits on both sides of (3.8) and interchange the order of limit and expectation. In cases N and P this is because of monotone convergence. Infinity is allowed as a possible limiting value. We obtain

$$\begin{aligned} F(\pi, x) &= c(x, u) + \beta E[F(\pi, x_1) \mid x_0 = x, u_0 = u] \\ &\geq c(x, u) + \beta E[F(x_1) \mid x_0 = x, u_0 = u] \\ &\geq \inf_u \{c(x, u) + \beta E[F(x_1) \mid x_0 = x, u_0 = u]\}. \end{aligned}$$

Minimizing the left hand side over π gives ' \geq '.

To prove ' \leq ', fix x and consider a policy π that having chosen u_0 and reached state x_1 then follows a policy π^1 which is suboptimal by less than ϵ from that point, i.e. $JF(\pi^1, x_1) \leq F(x_1) + \epsilon$. Note that such a policy must exist, by definition of F , although π^1 will depend on x_1 . We have

$$\begin{aligned} F(x) &\leq F(\pi, x) \\ &= c(x, u_0) + \beta E[F(\pi^1, x_1) \mid x_0 = x, u_0] \\ &\leq c(x, u_0) + \beta E[F(x_1) + \epsilon \mid x_0 = x, u_0] \\ &\leq c(x, u_0) + \beta E[F(x_1) \mid x_0 = x, u_0] + \beta\epsilon. \end{aligned}$$

Minimizing the right hand side over u_0 and recalling that ϵ is arbitrary gives ' \leq '. \square

3.5 Example: selling an asset

A speculator owns a rare collection of tulip bulbs and each day has an opportunity to sell it, which she may either accept or reject. The potential sale prices are independently and identically distributed with probability density function $g(x)$, $x \geq 0$. Each day there is a probability $1 - \beta$ that the market for tulip bulbs will collapse, making her bulb collection completely worthless. Find the policy that maximizes her expected return and express it as the unique root of an equation. Show that if $\beta > 1/2$, $g(x) = 2/x^3$, $x \geq 1$, then she should sell the first time the sale price is at least $\sqrt{\beta/(1 - \beta)}$.

Solution. There are only two states, depending on whether she has sold the collection or not. Let these be 0 and 1, respectively. The optimality equation is

$$\begin{aligned} F(1) &= \int_{y=0}^{\infty} \max[y, \beta F(1)] g(y) dy \\ &= \beta F(1) + \int_{y=0}^{\infty} \max[y - \beta F(1), 0] g(y) dy \\ &= \beta F(1) + \int_{y=\beta F(1)}^{\infty} [y - \beta F(1)] g(y) dy \end{aligned}$$

Hence

$$(1 - \beta)F(1) = \int_{y=\beta F(1)}^{\infty} [y - \beta F(1)] g(y) dy. \quad (3.9)$$

That this equation has a unique root, $F(1) = F^*$, follows from the fact that left and right hand sides are increasing and decreasing in $F(1)$, respectively. Thus she should sell when she can get at least βF^* . Her maximal reward is F^* .

Consider the case $g(y) = 2/y^3$, $y \geq 1$. The left hand side of (3.9) is less than the right hand side at $F(1) = 1$ provided $\beta > 1/2$. In this case the root is greater than 1 and we compute it as

$$(1 - \beta)F(1) = 2/\beta F(1) - \beta F(1)/[\beta F(1)]^2,$$

and thus $F^* = 1/\sqrt{\beta(1-\beta)}$ and $\beta F^* = \sqrt{\beta/(1-\beta)}$.

If $\beta \leq 1/2$ she should sell at any price.

Notice that discounting arises in this problem because at each stage there is a probability $1 - \beta$ that a ‘catastrophe’ will occur that brings things to a sudden end. This characterization of the way that discounting can arise is often quite useful.

What if past offers remain open? The state is now the best of the offers received and the dynamic programming equation is

$$\begin{aligned} F(x) &= \int_{y=0}^{\infty} \max[y, \beta F(\max(x, y))] g(y) dy \\ &= \int_{y=0}^x \max[y, \beta F(x)] g(y) dy + \int_{y=x}^{\infty} \max[y, \beta F(y)] g(y) dy \end{aligned}$$

However, the solution is exactly the same as before: sell at the first time an offer exceeds βF^* . Can you see why?

4 Positive Programming

Special theory for maximizing positive rewards. We see that there can be no optimal policy. However, if a given policy has a value function that satisfies the optimality equation then that policy is optimal. Value iteration algorithm.

4.1 Example: possible lack of an optimal policy.

Positive programming is about maximizing non-negative rewards, $r(x, u) \geq 0$, or minimizing non-positive costs, $c(x, u) \leq 0$. The following example shows that there may be no optimal policy.

Example 4.1. Suppose the possible states are the non-negative integers and in state x we have a choice of either moving to state $x + 1$ and receiving no reward, or moving to state 0, obtaining reward $1 - 1/x$, and then remaining in state 0 thereafter and obtaining no further reward. The optimality equations is

$$F(x) = \max\{1 - 1/x, F(x + 1)\} \quad x > 0.$$

Clearly $F(x) = 1$, $x > 0$, but the policy that chooses the maximizing action in the optimality equation always moves on to state $x + 1$ and hence has zero reward. Clearly, there is no policy that actually achieves a reward of 1.

4.2 Characterization of the optimal policy

The following theorem provides a necessary and sufficient condition for a policy to be optimal: namely, its value function must satisfy the optimality equation. This theorem also holds for the case of strict discounting and bounded costs.

Theorem 4.2. *Suppose D or P holds and π is a policy whose value function $F(\pi, x)$ satisfies the optimality equation*

$$F(\pi, x) = \sup_u \{r(x, u) + \beta E[F(\pi, x_1) \mid x_0 = x, u_0 = u]\}.$$

Then π is optimal.

Proof. Let π' be any policy and suppose it takes $u_t(x) = f_t(x)$. Since $F(\pi, x)$ satisfies the optimality equation,

$$F(\pi, x) \geq r(x, f_0(x)) + \beta E_{\pi'}[F(\pi, x_1) \mid x_0 = x, u_0 = f_0(x)].$$

By repeated substitution of this into itself, we find

$$F(\pi, x) \geq E_{\pi'} \left[\sum_{t=0}^{s-1} \beta^t r(x_t, u_t) \mid x_0 = x \right] + \beta^s E_{\pi'}[F(\pi, x_s) \mid x_0 = x]. \quad (4.1)$$

In case P we can drop the final term on the right hand side of (4.1) (because it is non-negative) and then let $s \rightarrow \infty$; in case D we can let $s \rightarrow \infty$ directly, observing that this term tends to zero. Either way, we have $F(\pi, x) \geq F(\pi', x)$. \square

4.3 Example: optimal gambling

A gambler has i pounds and wants to increase this to N . At each stage she can bet any whole number of pounds not exceeding her capital, say $j \leq i$. Either she wins, with probability p , and now has $i + j$ pounds, or she loses, with probability $q = 1 - p$, and has $i - j$ pounds. Let the state space be $\{0, 1, \dots, N\}$. The game stops upon reaching state 0 or N . The only non-zero reward is 1, upon reaching state N . Suppose $p \geq 1/2$. Prove that the timid strategy, of always betting only 1 pound, maximizes the probability of the gambler attaining N pounds.

Solution. The optimality equation is

$$F(i) = \max_{j: j \leq i} \{pF(i+j) + qF(i-j)\}.$$

To show that the timid strategy, say π , is optimal we need to find its value function, say $G(i) = F(\pi, x)$, and then show that it is a solution to the optimality equation. We have $G(i) = pG(i+1) + qG(i-1)$, with $G(0) = 0$, $G(N) = 1$. This recurrence gives

$$G(i) = \begin{cases} \frac{1 - (q/p)^i}{1 - (q/p)^N} & p > 1/2, \\ \frac{i}{N} & p = 1/2. \end{cases}$$

If $p = 1/2$, then $G(i) = i/N$ clearly satisfies the optimality equation. If $p > 1/2$ we simply have to verify that

$$G(i) = \frac{1 - (q/p)^i}{1 - (q/p)^N} = \max_{j: j \leq i} \left\{ p \left[\frac{1 - (q/p)^{i+j}}{1 - (q/p)^N} \right] + q \left[\frac{1 - (q/p)^{i-j}}{1 - (q/p)^N} \right] \right\}.$$

Let W_j be the expression inside $\{ \}$ on the right hand side. It is simple calculation to show that $W_{j+1} < W_j$ for all $j \geq 1$. Hence $j = 1$ maximizes the right hand side.

4.4 Value iteration

An important and practical method of computing F is **successive approximation** or **value iteration**. Starting with $F_0(x) = 0$, we can successively calculate, for $s = 1, 2, \dots$,

$$F_s(x) = \inf_u \{c(x, u) + \beta E[F_{s-1}(x_1) \mid x_0 = x, u_0 = u]\}.$$

So $F_s(x)$ is the infimal cost over s steps. Now let

$$F_\infty(x) = \lim_{s \rightarrow \infty} F_s(x) = \lim_{s \rightarrow \infty} \inf_{\pi} F_s(\pi, x) = \lim_{s \rightarrow \infty} \mathcal{L}^s(0). \quad (4.2)$$

This exists (by monotone convergence under N or P, or by the fact that under D the cost incurred after time s is vanishingly small.)

Notice that (4.2) reverses the order of $\lim_{s \rightarrow \infty}$ and \inf_{π} in (3.6). The following theorem states that we can interchange the order of these operations and that therefore $F_s(x) \rightarrow F(x)$. However, in case N we need an additional assumption:

F (finite actions): There are only finitely many possible values of u in each state.

Theorem 4.3. *Suppose that D or P holds, or N and F hold. Then $F_\infty(x) = F(x)$.*

Proof. First we prove ‘ \leq ’. Given any $\bar{\pi}$,

$$F_\infty(x) = \lim_{s \rightarrow \infty} F_s(x) = \lim_{s \rightarrow \infty} \inf_{\pi} F_s(\pi, x) \leq \lim_{s \rightarrow \infty} F_s(\bar{\pi}, x) = F(\bar{\pi}, x).$$

Taking the infimum over $\bar{\pi}$ gives $F_\infty(x) \leq F(x)$.

Now we prove ‘ \geq ’. In the positive case, $c(x, u) \leq 0$, so $F_s(x) \geq F(x)$. Now let $s \rightarrow \infty$. In the discounted case, with $|c(x, u)| < B$, imagine subtracting $B > 0$ from every cost. This reduces the infinite-horizon cost under any policy by exactly $B/(1 - \beta)$ and $F(x)$ and $F_\infty(x)$ also decrease by this amount. All costs are now negative, so the result we have just proved applies. [Alternatively, note that

$$F_s(x) - \beta^s B/(1 - \beta) \leq F(x) \leq F_s(x) + \beta^s B/(1 - \beta)$$

(can you see why?) and hence $\lim_{s \rightarrow \infty} F_s(x) = F(x)$.]

In the negative case,

$$\begin{aligned} F_\infty(x) &= \lim_{s \rightarrow \infty} \min_u \{c(x, u) + E[F_{s-1}(x_1) \mid x_0 = x, u_0 = u]\} \\ &= \min_u \{c(x, u) + \lim_{s \rightarrow \infty} E[F_{s-1}(x_1) \mid x_0 = x, u_0 = u]\} \\ &= \min_u \{c(x, u) + E[F_\infty(x_1) \mid x_0 = x, u_0 = u]\}, \end{aligned} \tag{4.3}$$

where the first equality follows because the minimum is over a finite number of terms and the second equality follows by Lebesgue monotone convergence (since $F_s(x)$ increases in s). Let π be the policy that chooses the minimizing action on the right hand side of (4.3). This implies, by substitution of (4.3) into itself, and using the fact that N implies $F_\infty \geq 0$,

$$\begin{aligned} F_\infty(x) &= E_\pi \left[\sum_{t=0}^{s-1} c(x_t, u_t) + F_\infty(x_s) \mid x_0 = x \right] \\ &\geq E_\pi \left[\sum_{t=0}^{s-1} c(x_t, u_t) \mid x_0 = x \right]. \end{aligned}$$

Letting $s \rightarrow \infty$ gives $F_\infty(x) \geq F(\pi, x) \geq F(x)$. □

4.5 Example: search for a moving object

Initially an object is equally likely to be in one of two boxes. If we search box 1 and the object is there we will discover it with probability 1, but if it is in box 2 and we search there then we will find it only with probability 1/2, and if we do not then the object moves to box 1 with probability 1/4. Suppose that we find the object on our N th search. Our aim is to maximize $E\beta^N$, where $0 < \beta < 1$.

Our state variable is p_t , the probability that the object is in box 1 given that we have not yet found it. If at time t we search box 1 and fail to find the object, then $p_{t+1} = 0$. On the other hand, if we search box 2,

$$p_{t+1} = a(p_t) = \frac{p_t + q_t(1/2)(1/4)}{p_t + (1/2)q_t} = \frac{1 + 7p_t}{8(1 - 0.5q_t)}.$$

The optimality equation is

$$F(p) = \max[p + q\beta F(0), 0.5q + (1 - 0.5q)\beta F(a(p))]. \quad (4.4)$$

By value iteration of

$$F_s(p) = \max[p + q\beta F_{s-1}(0), 0.5q + (1 - 0.5q)\beta F_{s-1}(a(p))]$$

we can prove that $F(p)$ is convex in p , by induction. The key fact is that $F_s(p)$ is always the maximum of a collection of linear function of p . We also use the fact that maximums of convex functions are convex.

The left hand side of (4.4) is linear in p taking values at $p = 0$ and $p = 1$ of $\beta F(0)$ and 1 respectively, and the right hand side takes values $0.5 + 0.5\beta F(.25)$ and $\beta F(0)$. Thus there is a unique p for which the left and right hand sides are equal, say p^* , and so an optimal policy is to search box 1 if and only if $p_t \geq p^*$.

4.6 Example: pharmaceutical trials

A doctor has two drugs available to treat a disease. One is well-established drug and is known to work for a given patient with probability p , independently of its success for other patients. The new drug is untested and has an unknown probability of success θ , which the doctor believes to be uniformly distributed over $[0, 1]$. He treats one patient per day and must choose which drug to use. Suppose he has observed s successes and f failures with the new drug. Let $F(s, f)$ be the maximal expected-discounted number of future patients who are successfully treated if he chooses between the drugs optimally from this point onwards. For example, if he uses only the established drug, the expected-discounted number of patients successfully treated is $p + \beta p + \beta^2 p + \dots = p/(1 - \beta)$. The posterior distribution of θ is

$$f(\theta | s, f) = \frac{(s + f + 1)!}{s!f!} \theta^s (1 - \theta)^f, \quad 0 \leq \theta \leq 1,$$

and the posterior mean is $\bar{\theta}(s, f) = (s + 1)/(s + f + 2)$. The optimality equation is

$$F(s, f) = \max \left[\frac{p}{1 - \beta}, \frac{s + 1}{s + f + 2} (1 + \beta F(s + 1, f)) + \frac{f + 1}{s + f + 2} \beta F(s, f + 1) \right].$$

Notice that after the first time that the doctor decides is not optimal to use the new drug it cannot be optimal for him to return to using it later, since his information about that drug cannot have changed while not using it.

It is not possible to give a closed-form expression for F , but we can find an approximate numerical solution. If $s + f$ is very large, say 300, then $\bar{\theta}(s, f) = (s + 1)/(s + f + 2)$ is a good approximation to θ . Thus we can take $F(s, f) \approx (1 - \beta)^{-1} \max[p, \bar{\theta}(s, f)]$, $s + f = 300$ and work backwards. For $\beta = 0.95$, one obtains the following table.

f	s	0	1	2	3	4	5
0		.7614	.8381	.8736	.8948	.9092	.9197
1		.5601	.6810	.7443	.7845	.8128	.8340
2		.4334	.5621	.6392	.6903	.7281	.7568
3		.3477	.4753	.5556	.6133	.6563	.6899
4		.2877	.4094	.4898	.5493	.5957	.6326

These numbers are the greatest values of p (the known success probability of the well-established drug) for which it is worth continuing with at least one more trial of the new drug. For example, suppose $p = 0.6$ and 6 trials with the new drug have given $s = f = 3$. Then since $p = 0.6 < 0.6133$ we should treat the next patient with the new drug. At this point the probability that the new drug will successfully treat the next patient is 0.5 and so the doctor will actually be treating that patient with the drug that is least likely to cure!

Here we see a tension going on between desires for **exploitation** and **exploration**. A **myopic policy** seeks only to maximize immediate reward. However, an optimal policy takes account of the possibility of gaining information that could lead to greater rewards being obtained later on. Notice that it is worth using the new drug at least once if $p < 0.7614$, even though at its first use the new drug will only be successful with probability 0.5. Of course as the discount factor β tends to 0 the optimal policy will look more and more like the myopic policy.

The above is an example of a **two-armed bandit problem** and a foretaste for Chapter 7 in which we will learn about the **multi-armed bandit problem** and how to optimally conduct trials amongst several alternative drugs with unknown success probabilities.

5 Negative Programming

The special theory of minimizing positive costs. We see that action that extremizes the right hand side of the optimality equation is an optimal policy. Stopping problems and their solution.

5.1 Example: a partially observed MDP

Example 5.1. Consider a similar problem to that of §4.5. A hidden object moves between two location according to a Markov chain with probability transition matrix $P = (p_{ij})$. A search in location i costs c_i , and if the object is there it is found with probability α_i . The aim is to minimize the expected cost of finding the object.

This is example of what is called a **partially observable Markov decision process** (POMDP). In a POMDP the decision-maker cannot directly observe the underlying state. Instead, he must maintain a probability distribution over the set of possible states, based on his observations, and the underlying MDP. This distribution is updated using the usual Bayesian calculations.

Solution. Let x_i be the probability that the object is in location i (where $x_1 + x_2 = 1$). Value iteration of the dynamic programming equation is via

$$F_s(x_1) = \min \left\{ c_1 + (1 - \alpha_1 x_1) F_{s-1} \left(\frac{(1 - \alpha_1) x_1 p_{11} + x_2 p_{21}}{1 - \alpha_1 x_1} \right), \right. \\ \left. c_2 + (1 - \alpha_2 x_2) F_{s-1} \left(\frac{(1 - \alpha_2) x_2 p_{21} + x_1 p_{11}}{1 - \alpha_2 x_2} \right) \right\}.$$

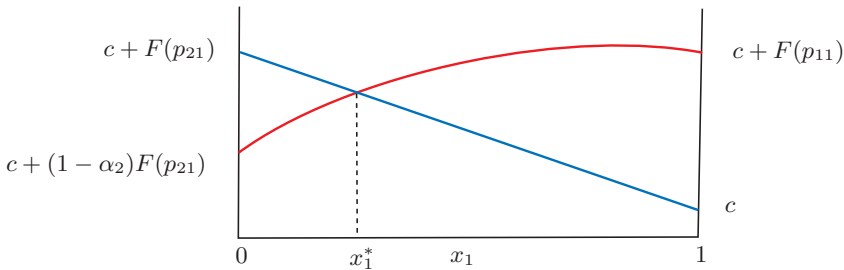
The arguments of $F_{s-1}(\cdot)$ are the posterior probabilities that the object in location 1, given that we have search location 1 (or 2) and not found it.

Now $F_0(x_1) = 0$, $F_1(x_1) = \min\{c_1, c_2\}$, $F_2(x)$ is the minimum of two linear functions of x_1 . If F_{s-1} is the minimum of some collection of linear functions of x_1 it follows that the same can be said of F_s . Thus, by induction, F_s is a concave function of x_1 .

By application of our theorem that $F_s \rightarrow F$ in the N and F case, we can deduce that the infinite horizon return function, F , is also a concave function. Notice that in the optimality equation for F (obtained by letting $s \rightarrow \infty$ in the equation above), the left hand term within the $\min\{\cdot, \cdot\}$ varies from $c_1 + F(p_{21})$ to $c_1 + (1 - \alpha_1)F(p_{11})$ as x_1 goes from 0 to 1. The right hand term varies from $c_2 + (1 - \alpha_2)F(p_{21})$ to $c_2 + F(p_{11})$ as x_1 goes from 0 to 1.

Consider the special case of $\alpha_1 = 1$ and $c_1 = c_2 = c$. Then the left hand term is the linear function $c + (1 - x_1)F(p_{21})$. This means we have the picture below, where the blue and red curves corresponds to the left and right hand terms, and intersect exactly once since the red curve is concave.

Thus the optimal policy can be characterized as “*search location 1 iff the probability that the object is in location 1 exceeds a threshold x_1^** ”.



The value of x_1^* depends on the parameters, α_i and p_{ij} . It is believed that the answer is of this form for any parameters, but this is still an unproved conjecture.

5.2 Stationary policies

A **Markov policy** is a policy that specifies the control at time t to be simply a function of the state and time. In the proof of Theorem 4.2 we used $u_t = f_t(x_t)$ to specify the control at time t . This is a convenient notation for a Markov policy, and we can write $\pi = (f_0, f_1, \dots)$ to denote such a policy. If in addition the policy does not depend on time and is non-randomizing in its choice of action then it is said to be a **deterministic stationary Markov policy**, and we write $\pi = (f, f, \dots) = f^\infty$.

For such a policy we might write

$$F_t(\pi, x) = c(x, f(x)) + E[F_{t+1}(\pi, x_1) \mid x_t = x, u_t = f(x)]$$

or $F_{t+1} = \mathcal{L}(f)F_{t+1}$, where $\mathcal{L}(f)$ is the operator having action

$$\mathcal{L}(f)\phi(x) = c(x, f(x)) + E[\phi(x_1) \mid x_0 = x, u_0 = f(x)].$$

5.3 Characterization of the optimal policy

Negative programming is about maximizing non-positive rewards, $r(x, u) \leq 0$, or minimizing non-negative costs, $c(x, u) \geq 0$. The following theorem gives a necessary and sufficient condition for a stationary policy to be optimal: namely, it must choose the optimal u on the right hand side of the optimality equation. Note that in this theorem we are requiring that the infimum over u is attained as a minimum over u (as would be the case if we make the finite actions assumptions, F).

Theorem 5.2. *Suppose D or N holds. Suppose $\pi = f^\infty$ is the stationary Markov policy such that*

$$f(x) = \arg \min_u [c(x, u) + \beta E[F(x_1) \mid x_0 = x, u_0 = u]].$$

Then $F(\pi, x) = F(x)$, and π is optimal.

(i.e. $u = f(x)$ is the value of u which minimizes the r.h.s. of the DP equation.)

Proof. The proof is really the same as the final part of proving Theorem 4.3. By substituting the optimality equation into itself and using the fact that π specifies the minimizing control at each stage,

$$F(x) = E_\pi \left[\sum_{t=0}^{s-1} \beta^t c(x_t, u_t) \middle| x_0 = x \right] + \beta^s E_\pi [F(x_s) | x_0 = x]. \quad (5.1)$$

In case N we can drop the final term on the right hand side of (5.1) (because it is non-negative) and then let $s \rightarrow \infty$; in case D we can let $s \rightarrow \infty$ directly, observing that this term tends to zero. Either way, we have $F(x) \geq F(\pi, x)$. \square

A corollary is that if assumption F holds then an optimal policy exists. Neither Theorem 5.2 or this corollary are true for positive programming (see Example 4.1).

5.4 Optimal stopping over a finite horizon

One way that the total-expected cost can be finite is if it is possible to enter a state from which no further costs are incurred. Suppose u has just two possible values: $u = 0$ (stop), and $u = 1$ (continue). Suppose there is a termination state, say 0, that is entered upon choosing the stopping action. Once this state is entered the system stays in that state and no further cost is incurred thereafter. We let $c(x, 0) = k(x)$ (stopping cost) and $c(x, 1) = c(x)$ (continuation cost). This defines a **stopping problem**.

Suppose that $F_s(x)$ denotes the minimum total cost when we are constrained to stop within the next s steps. This gives a finite-horizon problem with dynamic programming equation

$$F_s(x) = \min\{k(x), c(x) + E[F_{s-1}(x_1) | x_0 = x, u_0 = 1]\}, \quad (5.2)$$

with $F_0(x) = k(x)$, $c(0) = 0$.

Consider the set of states in which it is at least as good to stop now as to continue one more step and then stop:

$$S = \{x : k(x) \leq c(x) + E[k(x_1) | x_0 = x, u_0 = 1]\}.$$

Clearly, it cannot be optimal to stop if $x \notin S$, since in that case it would be strictly better to continue one more step and then stop. If S is closed then the following theorem gives us the form of the optimal policies for all finite-horizons.

Theorem 5.3. *Suppose S is closed (so that once the state enters S it remains in S .) Then an optimal policy for all finite horizons is: stop if and only if $x \in S$.*

Proof. The proof is by induction. If the horizon is $s = 1$, then obviously it is optimal to stop only if $x \in S$. Suppose the theorem is true for a horizon of $s - 1$. As above, if $x \notin S$ then it is better to continue for more one step and stop rather than stop in state x . If $x \in S$, then the fact that S is closed implies $x_1 \in S$ and so $F_{s-1}(x_1) = k(x_1)$. But then (5.2) gives $F_s(x) = k(x)$. So we should stop if $s \in S$. \square

The optimal policy is known as a **one-step look-ahead rule** (OSLA rule).

5.5 Example: optimal parking

A driver is looking for a parking space on the way to his destination. Each parking space is free with probability p independently of whether other parking spaces are free or not. The driver cannot observe whether a parking space is free until he reaches it. If he parks s spaces from the destination, he incurs cost s , $s = 0, 1, \dots$. If he passes the destination without having parked the cost is D . Show that an optimal policy is to park in the first free space that is no further than s^* from the destination, where s^* is the greatest integer s such that $(Dp + 1)q^s \geq 1$.

Solution. When the driver is s spaces from the destination it only matters whether the space is available ($x = 1$) or full ($x = 0$). The optimality equation gives

$$F_s(0) = qF_{s-1}(0) + pF_{s-1}(1),$$

$$F_s(1) = \min \begin{cases} s, & \text{(take available space)} \\ qF_{s-1}(0) + pF_{s-1}(1), & \text{(ignore available space)} \end{cases}$$

where $F_0(0) = D$, $F_0(1) = 0$.

Now we solve the problem using the idea of a OSLA rule. It is better to stop now (at a distance s from the destination) than to go on and take the first available space if s is in the stopping set

$$S = \{s : s \leq k(s-1)\}$$

where $k(s-1)$ is the expected cost if we take the first available space that is $s-1$ or closer. Now

$$k(s) = ps + qk(s-1),$$

with $k(0) = qD$. The general solution is of the form $k(s) = -q/p + s + cq^s$. So after substituting and using the boundary condition at $s = 0$, we have

$$k(s) = -\frac{q}{p} + s + \left(D + \frac{1}{p}\right)q^{s+1}, \quad s = 0, 1, \dots$$

So

$$S = \{s : (Dp + 1)q^s \geq 1\}.$$

This set is closed (since s decreases) and so by Theorem 5.3 this stopping set describes the optimal policy.

We might let D be the expected distance that that the driver must walk if he takes the first available space at the destination or further down the road. In this case, $D = 1 + qD$, so $D = 1/p$ and s^* is the greatest integer such that $2q^s \geq 1$.

6 Optimal Stopping Problems

More on stopping problems and their solution.

6.1 Bruss's odds algorithm

A doctor, using a special treatment, codes 1 for a successful treatment, 0 otherwise. He treats a sequence of n patients and wants to minimize any suffering, while achieving a success with every patient for whom that is possible. Stopping on the last 1 would achieve this objective, so he wishes to maximize the probability of this.

Solution. Suppose X_k is the code of the k th patient. Assume X_1, \dots, X_n are independent with $p_k = P(X_k = 1)$. Let $q_k = 1 - p_k$ and $r_k = p_k/q_k$. **Bruss's odds algorithm** sums the odds from the s th event to the last event (the n th)

$$R_s = r_s + \dots + r_n$$

and finds the greatest s , say s^* , for which $R_s \geq 1$. We claim that by stopping the first time that code 1 occurs amongst patients $\{s^*, s^* + 1, \dots, n\}$, the doctor maximizes probability of stopping on the last patient who can be successfully treated.

To prove this claim we just check optimality of a OSLA-rule. The stopping set is

$$\begin{aligned} S &= \{i : q_{i+1} \dots q_n > (p_{i+1}q_{i+2}q_{i+3} \dots q_n) + (q_{i+1}p_{i+2}q_{i+3} \dots q_n) \\ &\quad + \dots + (q_{i+1}q_{i+2}q_{i+3} \dots p_n)\} \\ &= \{i : 1 > r_{i+1} + r_{i+2} + \dots + r_n\} \\ &= \{s^*, s^* + 1, \dots, n\}. \end{aligned}$$

Clearly the stopping set is closed, so the OSLA-rule is optimal. The probability of stopping on the last 1 is $(q_{s^*} \dots q_n)(r_{s^*} + \dots + r_n)$ and (by solving a little optimization problem) this is always $\geq 1/e = 0.368$, provided $R_1 \geq 1$.

We can use the odds algorithm to re-solve the secretary problem. Code 1 when a candidate is better than all who have been seen previously. Our aim is to stop on the last candidate coded 1. We proved previously that X_1, \dots, X_h are independent and $P(X_t = 1) = 1/t$. So $r_i = (1/t)/(1 - 1/t) = 1/(t - 1)$. The algorithm tells us to ignore the first $s^* - 1$ candidates and the hire the first who is better than all we have seen previously, where s^* is the greatest integer s for which

$$\frac{1}{s-1} + \frac{1}{s} + \dots + \frac{1}{h-1} \geq 1 \quad \left(\equiv \text{the least } s \text{ for which } \frac{1}{s} + \dots + \frac{1}{h-1} \leq 1 \right).$$

We can also solve a 'groups' version of the secretary problem. Suppose we see h groups of candidates, of sizes n_1, \dots, n_h . We wish to stop with the group that contains the best of all the candidates. Then $p_1 = 1$, $p_2 = n_2/(n_1 + n_2), \dots, p_h = n_h/(n_1 + \dots + n_h)$. The odds algorithm tells us to stop if group i contains the best candidate so far and $i \geq s^*$, where s^* is the greatest integer such that

$$\frac{n_s}{\sum_{i=1}^{s-1} n_i} + \frac{n_{s+1}}{\sum_{i=1}^s n_i} + \dots + \frac{n_h}{\sum_{i=1}^{h-1} n_i} \geq 1.$$

6.2 Example: Stopping a random walk

indexstopping a random walk Suppose that x_t follows a random walk on $\{0, \dots, N\}$. At any time t we may stop the walk and take a positive reward $r(x_t)$. In states 0 and N we must stop. The aim is to maximize $Er(x_T)$.

Solution. The dynamic programming equation is

$$\begin{aligned} F(0) &= r(0), & F(N) &= r(N) \\ F(x) &= \max \left\{ r(x), \frac{1}{2}F(x-1) + \frac{1}{2}F(x+1) \right\}, & 0 < x < N. \end{aligned}$$

We see that

- (i) $F(x) \geq \frac{1}{2}F(x-1) + \frac{1}{2}F(x+1)$, so $F(x)$ is concave.
- (ii) Also $F(x) \geq r(x)$.

We say F is a **concave majorant** of r .

In fact, F can be characterized as the smallest concave majorant of r . For suppose that G is any other concave majorant of r .

Starting with $F_0 = 0$, we have $G \geq F_0$. So we can prove by induction that

$$\begin{aligned} F_s(x) &= \max \left\{ r(x), \frac{1}{2}F_{s-1}(x-1) + \frac{1}{2}F_{s-1}(x+1) \right\} \\ &\leq \max \left\{ r(x), \frac{1}{2}G(x-1) + \frac{1}{2}G(x+1) \right\} \\ &\leq \max \{ r(x), G(x) \} \\ &\leq G(x). \end{aligned}$$

Theorem 4.3 tells us that $F_s(x) \rightarrow F(x)$ as $s \rightarrow \infty$. Hence $F \leq G$.

A OSLA rule is not optimal here. The optimal rule is to stop iff $F(x) = r(x)$.

6.3 Optimal stopping over the infinite horizon

Consider now a general stopping problem over the infinite-horizon with $k(x), c(x)$ as previously, and with the aim of minimizing total expected cost. Let $F_s(x)$ be the infimal cost given that we are required to stop by the s th step. Let $F(x)$ be the infimal cost when all that is required is that we stop eventually. Since less cost can be incurred if we are allowed more time in which to stop, we have

$$F_s(x) \geq F_{s+1}(x) \geq F(x).$$

Thus by monotone convergence $F_s(x)$ tends to a limit, say $F_\infty(x)$, and $F_\infty(x) \geq F(x)$.

Example 6.1. Consider the problem of stopping a symmetric random walk on the integers, where $c(x) = 0$, $k(x) = \exp(-x)$. The policy of stopping immediately, say π , has $F(\pi, x) = \exp(-x)$, and since e^{-x} is a convex function this satisfies the infinite-horizon optimality equation,

$$F(x) = \min \{ \exp(-x), (1/2)F(x-1) + (1/2)F(x+1) \}.$$

However, π is not optimal. The random walk is recurrent, so we may wait until reaching as large an integer as we like before stopping; hence $F(x) = 0$. Thus we see two things:

- (i) It is possible that $F_\infty > F$. This is because $F_s(x) = e^{-x}$, but $F(x) = 0$.
- (ii) Theorem 4.2 is not true for negative programming. Policy π has $F(\pi, x) = e^{-x}$ and this satisfies the optimality equation. Yet π is not optimal.

Remark. In Theorem 4.3 we had $F_\infty = F$, but for that theorem we assumed $F_0(x) = k(x) = 0$ and $F_s(x)$ was the infimal cost possible over s steps, and thus $F_s \leq F_{s+1}$ (in the N case). However, Example 6.1 $k(x) > 0$ and $F_s(x)$ is the infimal cost amongst the set of policies that are required to stop within s steps. Now $F_s(x) \geq F_{s+1}(x)$.

The following lemma gives conditions under which the infimal finite-horizon cost does converge to the infimal infinite-horizon cost.

Lemma 6.2. *Suppose all costs are bounded as follows.*

$$(a) K = \sup_x k(x) < \infty \quad (b) C = \inf_x c(x) > 0. \quad (6.1)$$

Then $F_s(x) \rightarrow F(x)$ as $s \rightarrow \infty$.

Proof. Suppose π is an optimal policy for the infinite horizon problem and stops at the random time τ . It has expected cost of at least $(s+1)CP(\tau > s)$. However, since it would be possible to stop at time 0 the cost is also no more than K , so

$$(s+1)CP(\tau > s) \leq F(x) \leq K.$$

In the s -horizon problem we could follow π , but stop at time s if $\tau > s$. This implies

$$F(x) \leq F_s(x) \leq F(x) + KP(\tau > s) \leq F(x) + \frac{K^2}{(s+1)C}.$$

By letting $s \rightarrow \infty$, we have $F_\infty(x) = F(x)$. □

Note that the problem posed here is identical to one in which we pay K at the start and receive a terminal reward $r(x) = K - k(x)$.

Theorem 6.3. *Suppose S is closed and (6.1) holds. Then an optimal policy for the infinite horizon is: stop if and only if $x \in S$.*

Proof. By Theorem 5.3 we have for all finite s ,

$$F_s(x) = \begin{cases} k(x) & x \in S, \\ < k(x) & x \notin S. \end{cases}$$

Lemma 6.2 gives $F(x) = F_\infty(x)$. □

6.4 Sequential Probability Ratio Test

A statistician wishes to decide between two hypotheses, $H_0 : f = f_0$ and $H_1 : f = f_1$ on the basis of i.i.d. observations drawn from a distribution with density f . *Ex ante* he believes the probability that H_i is true is p_i (where $p_0 + p_1 = 1$). Suppose that he has the sample $x = (x_1, \dots, x_n)$. The posterior probabilities are in the likelihood ratio

$$\ell_n(x) = \frac{f_1(x_1) \cdots f_1(x_n) p_1}{f_0(x_1) \cdots f_0(x_n) p_0}.$$

Suppose it costs γ to make an observation. Stopping and declaring H_i true results in a cost c_i if wrong. This leads to the optimality equation for minimizing expected cost

$$F(\ell) = \min \left\{ c_0 \frac{\ell}{1+\ell}, c_1 \frac{1}{1+\ell}, \right. \\ \left. \gamma + \frac{\ell}{1+\ell} \int F(\ell f_1(y)/f_0(y)) f_1(y) dy + \frac{1}{1+\ell} \int F(\ell f_1(y)/f_0(y)) f_0(y) dy \right\}$$

Taking $H(\ell) = (1 + \ell)F(\ell)$, the optimality equation can be rewritten as

$$H(\ell) = \min \left\{ c_0 \ell, c_1, (1 + \ell)\gamma + \int H(\ell f_1(y)/f_0(y)) f_0(y) dy \right\}.$$

We have a very similar problem to that of searching for a moving object. The state is ℓ_n . We can stop (in two ways) or continue by paying for another observation, in which case the state makes a random jump to $\ell_{n+1} = \ell_n f_1(x)/f_0(x)$, where x is a sample from f_0 . We can show that $H(\cdot)$ is concave in ℓ , and that therefore the optimal policy can be described by two numbers, $a_0^* < a_1^*$: *If $\ell_n \leq a_0^*$, stop and declare H_0 true; If $\ell_n \geq a_1^*$, stop and declare H_1 true; otherwise take another observation.*

6.5 Bandit processes

A **bandit process** is a special type of MDP in which there are just two possible actions: $u = 0$ (freeze) or $u = 1$ (continue). The control $u = 0$ produces no reward and the state does not change (hence the term ‘freeze’). Under $u = 1$ we obtain a reward $r(x_t)$ and the state changes, to x_{t+1} , according to the Markov dynamics $P(x_{t+1} | x_t, u_t = 1)$.

A **simple family of alternative bandit processes** (SFABP) is a collection of n such bandit processes. At each time $t = 0, 1, \dots$ we must select exactly one bandit to receive continuation, while all others are frozen.

This is a rich modelling framework. With it we can model questions like this:

- Which of n drugs should we give to the next patient?
- Which of n jobs should we work on next?
- When of n oil fields should we explore next?

6.6 Example: Two-armed bandit

Consider a family of two alternative bandit processes. Bandit process B_1 is trivial: it stays in the same state and always produces known reward λ at each step that it is continued. Bandit process B_2 is nontrivial. It starts in state $x(0)$ and evolves as a Markov chain when it is continued and produces a state-dependent reward. The state of B_2 is what's important. Starting B_2 in state $x(0) = x$ we have the optimality equation

$$\begin{aligned} F(x) &= \max \left\{ \frac{\lambda}{1-\beta}, r(x) + \beta \sum_y P(x, y) F(y) \right\} \\ &= \max \left\{ \frac{\lambda}{1-\beta}, \sup_{\tau > 0} E \left[\sum_{t=0}^{\tau-1} \beta^t r(x(t)) + \beta^\tau \frac{\lambda}{1-\beta} \mid x(0) = x \right] \right\}. \end{aligned}$$

The left hand choice within $\max\{\cdot, \cdot\}$ corresponds to continuing B_1 . The right hand choice corresponds to continuing B_2 for at least one step and then switching to B_1 a some later step, τ . Notice that once we switch to B_1 we will never wish switch back to B_2 because things remain the same as when we first switched away from B_2 .

We are to choose the **stopping time** τ so as to optimally switch from continuing B_2 to continuing B_1 . Because the two terms within the $\max\{\cdot, \cdot\}$ are both increasing in λ , and are linear and convex, respectively, there is a unique λ , say λ^* , for which they are equal. This is

$$\lambda^* = \sup \left\{ \lambda : \frac{\lambda}{1-\beta} \leq \sup_{\tau > 0} E \left[\sum_{t=0}^{\tau-1} \beta^t r(x(t)) + \beta^\tau \frac{\lambda}{1-\beta} \mid x(0) = x \right] \right\}. \quad (6.2)$$

Of course this λ^* depends on $x(0)$. We denote its value as $G(x)$. After a little algebra

$$G(x) = \sup_{\tau > 0} \frac{E \left[\sum_{t=0}^{\tau-1} \beta^t r(x(t)) \mid x(0) = x \right]}{E \left[\sum_{t=0}^{\tau-1} \beta^t \mid x(0) = x \right]}.$$

G is called a **Gittins index**.

So we now have the complete solution to the two-armed bandit problem. *If $G(x(0)) \leq G$ then it is optimal to continue B_1 forever. If $G(x(0)) > G$ then it is optimal to continue B_2 until the first time τ at which $G(x(\tau)) \leq G$.*

6.7 Example: prospecting

We run a business that returns R_0 per day. We are considering making a change that might produce a better return of R_1 per day. Initially we know only that R_1 is distributed $U[0, 1]$. To trial this method for one day will cost c_1 , and at the end of this

day we will know R_1 . The Gittins index for the new method is G_1 such that

$$\begin{aligned} \frac{G_1}{1-\beta} &= -c_1 + E[R_1] + \frac{\beta}{1-\beta} E \max \{G_1, R_1\} \\ &= -c_1 + 1/2 + \frac{\beta}{1-\beta} \left[\int_0^{G_1} G_1 dr + \int_{G_1}^1 r dr \right]. \end{aligned}$$

For $\beta = 0.9$ and $c_1 = 1$ this gives $G_1 = 0.5232$. So it is worth conducting the trial only if $R_0 = G_0 \leq 0.5232$.

Now suppose that there is also a second method we might try. It produces reward R_2 , which is *ex ante* distributed $U[0, 2]$ and costs $c_2 = 3$ to trial. Its Gittins index is

$$\begin{aligned} \frac{G_2}{1-\beta} &= -c_2 + E[R_2] + \frac{\beta}{1-\beta} E \max \{G_2, R_2\} \\ &= -c_2 + 1 + \frac{\beta}{1-\beta} \left[\int_0^{G_2} G_2 \frac{1}{2} dr + \int_{G_2}^2 r \frac{1}{2} dr \right]. \end{aligned}$$

For $c_2 = 3$ this gives $G_2 = 0.8705$. Recall $G_0 = R_0$. Suppose $G_0 < G_1$. Since $G_2 > G_1 > G_0$ we might conjecture the following is optimal.

We start by trialing method 2. If $R_2 > G_1 = 0.5232$ stop and use method 2 thereafter. Otherwise, we trial method 1 and then, having learned all of R_0, R_1, R_2 , we pick the method that produces the greatest return and use that method thereafter.

However, this solution only a conjecture. We will prove it is optimal using the Gittins index theorem.

7 Bandit Processes and the Gittins Index

The multi-armed bandit problem. Bandit processes. Gittins index theorem.

7.1 Index policies

Recall the **single machine scheduling** example in §3.2 in which n jobs are to be processed successively on one machine. Job i has a known processing time t_i , assumed to be a positive integer. On completion of job i a positive reward r_i is obtained. We used an interchange argument to show that total discounted reward obtained from the n jobs is maximized by the **index policy** of always processing the uncompleted job of greatest index, computed as $r_i \beta^{t_i} (1 - \beta) / (1 - \beta^{t_i})$.

Notice that if we were allowed to interrupt the processing a job before finishing, so as to carry out some processing on a different job, this would be against the advice of the index policy. For the index, $r_i \beta^{t_i} (1 - \beta) / (1 - \beta^{t_i})$, increases as t_i decreases.

7.2 Multi-armed bandit problem

A **multi-armed bandit** is a slot-machine with multiple arms. The arms differ in the distributions of rewards that they pay out when pulled. An important special case is when arm i is a so-called **Bernoulli bandit**, with parameter p_i . We have already met this as the drug-testing model in §4.6. Such an arm pays $\mathcal{L}1$ with probability p_i , and $\mathcal{L}0$ with probability $1 - p_i$; this happens independently each time the arm is pulled. If there are n such arms, and a gambler knows the true values of p_1, \dots, p_n , then obviously he maximizes his expected reward by always pulling the arm of maximum p_i . However, if he does not know these values, then he must choose each successive arm on the basis of the information he has obtained by playing, updated in a Bayesian manner on the basis of observing the rewards he has obtained on previous pulls. The aim in the **multi-armed bandit problem** (MABP) is to maximize the expected total discounted reward.

More generally, we consider a problem of controlling the evolution of n independent reward-producing Markov processes decision processes. The action space of each process contains just two controls, which cause the process to be either ‘continued’ or ‘frozen’. At each instant (in discrete time) exactly one of these so-called **bandit processes** is *continued* (and reward from it obtained), while all the other bandit processes are *frozen*. The continued process can change state; but frozen processes do not change state. Reward is accrued only from the bandit process that is continued. This creates what is termed a **simple family of alternative bandit processes** (SFABP). The word ‘simple’ means that all the n bandit processes are available at all times.

Let $x(t) = (x_1(t), \dots, x_n(t))$ be the states of the n bandits. Let i_t denote the bandit process that is continued at time t under some policy π . In the language of Markov decision problems, we wish to find the value function:

$$F(x) = \sup_{\pi} E \left[\sum_{t=0}^{\infty} r_{i_t}(x_{i_t}(t)) \beta^t \mid x(0) = x \right],$$

where the supremum is taken over all policies π that are realizable (or non-anticipatory), in the sense that i_t depends only on the problem data and $x(t)$, not on any information which only becomes known only after time t .

Setup in this way, we have an infinite-horizon discounted-reward Markov decision problem. It therefore has a deterministic stationary Markov optimal policy. Its dynamic programming is

$$F(x) = \max_{i:i \in \{1, \dots, n\}} \left\{ r_i(x_i) + \beta \sum_{y \in E_i} P_i(x_i, y) F(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n) \right\}. \quad (7.1)$$

7.3 Gittins index theorem

Remarkably, the problem posed by a SFABP (or a MABP) can be solved by an **index policy**. That is, we can compute a number (called an index), separately for each bandit process, such that the optimal policy is always to continue the bandit process having the currently greatest index.

Theorem 7.1 (Gittins Index Theorem). *The problem posed by a SFABP, as setup above, is solved by always continuing the process having the greatest **Gittins index**, which is defined for bandit process i as*

$$G_i(x_i) = \sup_{\tau > 0} \frac{E \left[\sum_{t=0}^{\tau-1} \beta^t r_i(x_i(t)) \mid x_i(0) = x_i \right]}{E \left[\sum_{t=0}^{\tau-1} \beta^t \mid x_i(0) = x_i \right]}, \quad (7.2)$$

where τ is a stopping time constrained to take a value in the set $\{1, 2, \dots\}$.

The Index Theorem above is due to Gittins and Jones, who had obtained it by 1970, and presented it in 1972. The solution of the MABP impressed many experts as surprising and beautiful. Peter Whittle describes a colleague of high repute, asking another colleague ‘*What would you say if you were told that the multi-armed bandit problem had been solved?*’ The reply was ‘*Sir, the multi-armed bandit problem is not of such a nature that it can be solved*’.

The optimal **stopping time** τ in (7.2) is $\tau = \min\{t : G_i(x_i(t)) \leq G_i(x_i(0)), \tau > 0\}$, that is, τ is the first time at which the process reaches a state whose Gittins index is no greater than Gittins index at $x_i(0)$.

Examining (7.2), we see that the Gittins index is the maximal possible quotient of a numerator that is ‘expected total discounted *reward* over τ steps’, and denominator that is ‘expected total discounted *time* over τ steps’, where τ is at least 1 step. Notice that the Gittins index can be computed for all states of B_i as a function only of the data $r_i(\cdot)$ and $P_i(\cdot, \cdot)$. That is, it can be computed without knowing anything about the other bandit processes.

In the single machine scheduling example of §7.1, the optimal stopping time on the right hand side of (7.2) is $\tau = t_i$, the numerator is $r_i \beta^{t_i}$ and the denominator is $1 + \beta + \dots + \beta^{t_i-1} = (1 - \beta^{t_i}) / (1 - \beta)$. Thus, $G_i = r_i \beta^{t_i} (1 - \beta) / (1 - \beta^{t_i})$. Note that $G_i \rightarrow r_i / t_i$ as $\beta \rightarrow 1$.

7.4 Calibration

An alternative characterization of $G_i(x_i)$ is the one in (6.2)

$$G_i(x_i) = \sup \left\{ \lambda : \frac{\lambda}{1-\beta} \leq \sup_{\tau > 0} E \left[\sum_{t=0}^{\tau-1} \beta^t r_i(x_i(t)) + \beta^\tau \frac{\lambda}{1-\beta} \mid x_i(0) = x_i \right] \right\}. \quad (7.3)$$

That is, we consider a simple family of two bandit processes: bandit process B_i and a **calibrating bandit process**, say Λ , which pays out a known reward λ at each step it is continued. The Gittins index of B_i is the value of λ for which we are indifferent as to which of B_i and Λ to continue initially. Notice that once we decide to switch from continuing B_i to continuing Λ , at time τ , then information about B_i does not change and so it must be optimal to stick with continuing Λ ever after.

7.5 Proof of the Gittins index theorem

Various proofs have been given of the index theorem, all of which are useful in developing insight about this remarkable result. The following one is due to Weber (1992).

Proof of Theorem 7.1. We start by considering a problem in which only bandit process B_i is available. Let us define the **fair charge**, $\gamma_i(x_i)$, as the maximum amount that an agent would be willing to pay per step if he must continue B_i for one more step, and then stop whenever he likes thereafter. This is

$$\gamma_i(x_i) = \sup \left\{ \lambda : 0 \leq \sup_{\tau > 0} E \left[\sum_{t=0}^{\tau-1} \beta^t (r_i(x_i(t)) - \lambda) \mid x_i(0) = x_i \right] \right\}. \quad (7.4)$$

Notice that (7.3) and (7.4) are equivalent and so $\gamma_i(x_i) = G_i(x_i)$. Notice also that the time τ will be the first time that $G_i(x_i(\tau)) < G_i(x_i(0))$.

We next define the **prevailing charge** for B_i at time t as $g_i(t) = \min_{s \leq t} \gamma_i(x_i(s))$. So $g_i(t)$ actually depends on $x_i(0), \dots, x_i(t)$ (which we omit from its argument for convenience). Note that $g_i(t)$ is a nonincreasing function of t and its value depends only on the states through which bandit i evolves. The proof of the Index Theorem is completed by verifying the following facts, each of which is almost obvious.

- (i) Suppose that in the problem with n available bandit processes, B_1, \dots, B_n , the agent not only collects rewards, but also pays the prevailing charge of whatever bandit that he chooses to continue at each step. Then he cannot do better than just break even (i.e. expected value of rewards minus prevailing charges is 0).

This is because he could only make a strictly positive profit (in expected value) if this were to happen for at least one bandit. Yet the prevailing charge has been defined in such a way that he can only just break even.

- (ii) If he always continues the bandit of greatest prevailing charge then he will interleave the n nonincreasing sequences of prevailing charges into a single nonincreasing sequence of prevailing charges and so maximize their discounted sum.

- (iii) Using this strategy he also just breaks even; so this strategy, (of always continuing the bandit with the greatest $g_i(x_i)$), must also maximize the expected discounted sum of the rewards can be obtained from this SFABP. \square

7.6 Example: Weitzman's problem

'Pandora' has n boxes, each of which contains an unknown prize. *Ex ante* the prize in box i has a value with probability distribution function F_i . She can learn the value of the prize by opening box i , which costs her c_i to do. At any stage she may stop and take as her reward the maximum of the prizes she has found. She wishes to maximize the expected value of the prize she takes, minus the costs of opening boxes.

Solution. This problem is the 'prospecting' problem we already considered in §6.7. It can be modelled in terms of a SFABP. Box i is associated with a bandit process B_i , which starts in state 0. The first time it is continued there is a cost c_i , and the state becomes x_i , chosen by the distribution F_i . At all subsequent times that it is continued the reward is $r(x_i) = (1 - \beta)x_i$, and the state remains x_i . We wish to maximize the expected value of

$$-\sum_{t=1}^{\tau} \beta^{t-1} c_{i_t} + \max\{r(x_{i_1}), \dots, r(x_{i_\tau})\} \sum_{t=\tau}^{\infty} \beta^t$$

where we open boxes i_1, \dots, i_τ and then take the best prize thereafter. In the limit as $\beta \rightarrow 1$ this objective corresponds to that of Weitzman's problem, namely,

$$-\sum_{t=1}^{\tau} c_{i_t} + \max\{x_{i_1}, \dots, x_{i_\tau}\}$$

and so we can find the solution using the Gittins index theorem.

The Gittins index of an opened box is $r(x_i)$. The index of an unopened box i is the solution to

$$\frac{G_i}{1 - \beta} = -c_i + \frac{\beta}{1 - \beta} E \max\{r(x_i), G_i\}$$

or, by setting $g_i = G/(1 - \beta)$, and letting $\beta \rightarrow 1$, we get an index that is the solution of $g_i = -c_i + E \max\{x_i, g_i\}$.

For example, if F_i is a two point distribution with $x_i = 0$ or $x_i = r_i$, with probabilities $1 - p_i$ and p_i , then $g_i = -c_i + (1 - p_i)g_i + p_i r_i \implies g_i = r_i - c_i/p_i$.

Pandora's optimal strategy is thus: *Open boxes in decreasing order of g_i until first reaching a point that a revealed prize is greater than all g_i of unopened boxes.*

8 Applications of Bandit Processes

We consider some applications and generalizations to tax problems, job scheduling and branching bandits.

8.1 Forward induction policies

If we put $\tau = 1$ on the right hand side of (7.2) then it evaluates to $Er_i(x_i(t))$. If we use this as an index for choosing between projects, we have a **myopic policy** or **one-step-look-ahead policy**. The Gittins index policy generalizes the idea of a one-step-look-ahead policy, since it looks-ahead by some optimal time τ , so as to maximize, on the right hand side of (7.2), a measure of the rate at which reward can be accrued. This defines a so-called **forward induction policy**.

8.2 Example: playing golf with more than one ball

A golfer is playing with n balls. The balls are at positions x_1, \dots, x_n . If he plays ball i it will next land at location y , where $P(x_i, y)$ is known. He wishes to minimize the expected number of shots required to get one ball in the hole (location 0).

Solution. To represent this as a SFABP we shall set rewards and costs all 0, except that a reward R is obtained by continuing a bandit that is in state 0. So if some bandit reaches state 0, say with the golfer's t th shot, he will continue to play it there, obtaining reward $(\beta^t + \beta^{t+1} + \dots)R$. Suppose the golfer pays at the start a 'green fee' of $R/(1 - \beta)$. Then he will be trying to maximize

$$-\frac{R}{1 - \beta} + (\beta^t + \beta^{t+1} + \dots)R = -(1 + \beta + \dots + \beta^{t-1})R$$

which tends to $-tR$ as $\beta \rightarrow 1$. So he will be minimizing the expected number of shots needed to sink a ball. Locations are ordered by their Gittins indices. Location 0 has the greatest index, namely $G_0 = R/(1 - \beta)$. The golfer should always play the ball in location having greatest Gittins index.

Remark. We can reprise a proof of the index theorem, but working it for this golfing problem. Suppose the golfer is playing with just one ball, which is in location x_i . The golfer faces a cost of 1 for each shot he takes until the ball is sunk. So to motivate him to play, we offer a prize $g(x_i)$ which he wins if he plays at least one more shot and eventually sinks the ball. However, he may still quit if subsequently the ball lands in a bad place and the offered prize is no longer sufficiently motivating. If, however, that ever happens, we will increase the offered prize, so that it again becomes just advantageous for him to keep playing. This defines an nondecreasing sequence of offered prizes for ball i . Notice that they are defined independently of the other balls.

Now he plays with n balls. To each ball we attach an offered prize, just as above. It is a function of the ball's location, just as if he were playing only with that ball.

The key idea is that with these offered prizes the golfer can keep playing until some ball is sunk, and he will just break even. He is guaranteed to collect the least prize at the time a ball is finally sunk if he follows the policy of always playing the ball for which the least prize is offered. But the prizes were invented to make the game is ‘just fair’, and so in minimizing the value of the prize obtained when a ball is sunk this policy must also minimize the expected number of shots required until a ball is sunk. The prize $g(x_i)$ is of course the Gittins index for location x_i .

8.3 Target processes

The ‘problem above is different to our original set up of a SFABP problem. The golfing problem ends once one of the balls reaches the hole, and there is no discounting. The first issue we have modelled by allowing play to continue forever but making sure that it is optimal to keep playing the ball that is already in the hole, gaining R each time.

To introduce discounting we might take $P(x, 0) = 1 - \beta$ for all x . We also might easily generalize to there being a cost $c(x)$ for playing a ball in location x . The problem is now one of seeking to minimize

$$\begin{aligned} E \left[\sum_{t=0}^{\infty} \beta^t c(x_{i_t}(t)) - (1 - \beta) \frac{R}{1 - \beta} - \beta(1 - \beta) \frac{R}{1 - \beta} - \beta^2(1 - \beta) \frac{R}{1 - \beta} - \dots \right] \\ = E \left[\sum_{t=0}^{\tau-1} \beta^t c(x_{i_t}(t)) \right] - \frac{R}{1 - \beta}, \end{aligned}$$

where i_t is the ball played at time t and $x_{i_t}(t)$ is its state. Ignoring the term of $-R/(1-\beta)$, the objective function now looks exactly like one with which we are familiar.

The golfing problem is an example of a so-called **target process**. The aim is to control a SFABP in such a way as to minimize the expected cost incurred until such time as one of the bandits achieves some objective (such as landing in a certain state). For example, we might be viewing apartments in two neighborhoods, seeing their ‘values’ as two i.i.d. sequences of variables $X_1, X_2, \dots \sim F$, and $Y_1, Y_2, \dots \sim G$, but with F and G unknown. We wish to minimize the expected number of samples required to find one (a X_t or Y_t) that takes a value $\geq T$, for some given target T .

8.4 Bandit superprocesses

Suppose that a ball in location x can be played with a choice of golf clubs. If club $a \in A(x)$, is used then $x \rightarrow y$ with probability $P_a(x, y)$. Now the golfer must choose, not only which ball to play, but with which club to play it. Under a condition, an index policy is again optimal. He should play the ball with least prevailing prize, choosing the club from A that is optimal if that ball were the only ball present.

However, the condition for this policy to be optimal quite demanding. Recall that we defined $g(x)$ as a prize that the golfer could obtain by playing the ball for at least one more shot. The required condition is that whatever size of prize we offer, the golfer will optimally choose the same club.

8.5 Example: single machine stochastic scheduling

A collection of n jobs is to be processed on a single machine. They have unknown processing times, but we know how much processing each has already received, say x_1, \dots, x_n , units respectively. The probability that job i of age x_i will complete when it is next serviced is $h(x_i)$, the **hazard rate**. We wish to maximize the expected value of $\sum_i r_i \beta^{t_i}$, where t_i is the time at which job i completes, ($0 < \beta < 1$).

Solution. The jobs can be viewed as bandit processes. In computing the index we think about processing the job τ times, or until it completes, whichever comes first.

$$G(x_i) = \sup_{\tau > 0} \frac{\sum_{t=0}^{\tau-1} \beta^t r_i h(x_i + t) \prod_{s=0}^{t-1} (1 - h(x_i + s))}{\sum_{t=0}^{\tau-1} \beta^t \prod_{s=0}^{t-1} (1 - h(x_i + s))}.$$

There are two special cases for which it is easy to see what to do.

- $h(x)$ decreasing. Here $\tau = 1$ and $G(x_i) = r_i h(x_i)$.

We should always work on the job with the greatest index. Notice that this policy will be **preemptive**. We may wish to leave a job before it is finished.

- $h(x)$ increasing. Here $\tau = \infty$.

We should always work on the job with the greatest index and then continue processing it until it is complete. This policy will be **nonpreemptive**. Notice that as $\beta \rightarrow 1$, G tends to the r_i/ET_i , where T_i is the remaining processing time of job i .

Let C_i be the time that job i is completed. Notice also that setting $\alpha = 1 - \beta$,

$$\sum_i r_i \beta^{C_i} = \sum_i r_i (1 - \alpha)^{C_i} = \sum_i r_i - \alpha \sum_i r_i C_i + o(\alpha).$$

So in the limit $\beta \rightarrow 1$ we are solving a problem of minimizing $\sum_i r_i EC_i$, which is the expected value of a weighted sum of the completion times. This is known as the **weighted flow time**.

It might be more natural to replace r_i with c_i since we are seeking to minimize a cost. In the case that hazard rates are constant, but differ from job to job, say $h_i = \mu_i$ then jobs have processing times that are geometrically distributed and $ET_i = 1/\mu_i$ and the prescription is to process jobs in decreasing order of the index $c_i/ET_i = c_i \mu_i$. This is the famous so-called **$c\mu$ -rule**.

8.6 Calculation of the Gittins index

We can compute the Gittins indices in the following way. The input is the data of $r_i(\cdot)$ and $P_i(\cdot, \cdot)$. If the state space of B_i is finite, say $E_i = \{1, \dots, k_i\}$, then the Gittins indices can be computed in an iterative fashion. First we find the state of greatest index, say 1 such that $1 = \arg \max_j r_i(j)$. Having found this state we find the state of second-greatest index. If this is state j , then $G_i(j)$ is computed in (7.2) by taking τ to

be the first time that the state is not 1. This means that the second-best state is the state j which maximizes

$$\frac{E[r_i(j) + \beta r_i(1) + \cdots + \beta^{\tau-1} r_i(1)]}{E[1 + \beta + \cdots + \beta^{\tau-1}]},$$

where τ is the time at which, having started at $x_i(0) = j$, we have $x_i(\tau) \neq 1$. One can continue in this manner, successively finding states, and their Gittins indices, in decreasing order of their indices. If B_i moves on a finite state space of size k_i then its Gittins indices (one for each of the k_i states) can be computed in time $O(k_i^3)$.

If the state space of a bandit process is infinite, as in the case of the Bernoulli bandit, there may be no finite calculation by which to determine the Gittins indices for all states. In this circumstance, we can approximate the Gittins index using something like the value iteration algorithm. Essentially, one solves a problem of maximizing right hand side of (7.2), subject to $\tau \leq N$, where N is large.

8.7 Branching bandits

Consider a queue on n jobs. Job i has a deterministic processing time t_i and a reward r_i is obtained on completion. If we process the jobs in order $1, 2, \dots, n$ the completion time of job i is $C_i = t_1 + \cdots + t_i$ and the discounted sum of rewards is

$$r_1 \beta^{C_1} + r_2 \beta^{C_2} + \cdots + r_n \beta^{C_n}. \quad (8.1)$$

Notice that if we put $\beta = 1$ the above is independent of the order in which the processing is scheduled so the problem of optimal scheduling is vacuous. However, we can divide (8.1) by $1 - \beta$, subtract it from $(r_1 + \cdots + r_n)/(1 - \beta)$, and then let $\beta \rightarrow 1$ we get

$$r_1 C_1 + r_2 C_2 + \cdots + r_n C_n$$

which is a r_i -weighted sum of the completion times. It is now a non-vacuous problem to minimize this, which is done by always processing the job for which the index $r_i/t_i = \lim_{\beta \rightarrow 1} (1 - \beta)r_i \beta^{t_i} / (1 - \beta^{t_i})$ is greatest. This reprises our finding at the end of §8.5. We call this a **tax problem** (because the uncompleted jobs are taxed.)

Now we consider a problem in which the number of bandits is not fixed. Suppose the bandits are jobs of k different types and we start with n_i jobs of type i . Processing a job of type i takes time t_i and at the end of this processing a random number of new jobs arrive. Let's suppose that the number of newly arriving type j is distributed as a Poisson random variable with mean $\lambda_j t_i$. We pay a holding cost of r_i per job of type i that is in the system and wish to minimize the average holding cost. This is an example of what is called a **branching bandit** and the Gittins index theorem can be shown to hold.

Suppose that $r_1/t_1 > \cdots > r_n/t_n$. The Gittins indices can be computed in the manner of §8.6. It is easy to see that as long as there is at least one type 1 job in the system then it should be processed. Suppose we have emptied the system of type 1

jobs and are trying to discover which job type is of next greatest priority. If we process a job of type i then upon its completion we should next process all jobs of type 1 that have newly arrived until such point that the system is again cleared of all type 1 jobs. The expected number of jobs of type 1 that we will need to clear is proportional to t_i , say αt_i , the Gittins index is therefore (in the undiscounted case)

$$G_i = \frac{r_i + (\alpha t_i)r_1}{t_i + (\alpha t_i)t_1} = \frac{r_i/t_i + \alpha r_1}{1 + \alpha t_1}.$$

This is an increasing function of r_i/t_i , and so it is greatest for $i = 2$. So job type 2 has second greatest priority. Continuing in this way we see that the average waiting cost in a $M/G/1$ queue is minimized by always processing the job of greatest r_i/t_i . The Gittins indices are different to those in the problem with no arrivals, but they are ordered the same.

8.8 Example: Searching for a single object

An object is hidden in one of n boxes. *Ex ante*, it is in box i with probability p_i (where these sum to 1). The probability that a search in box i finds the object if it is there is q_i . We wish to minimize the expected number of searches needed to find the object.

Solution. This is not obviously a bandit problem. Looking box i and not finding the object changes all the posterior probabilities (indeed $p_i \rightarrow (1 - q_i)p_i/(1 - q_i p_i)$ and $p_j \rightarrow p_j/(1 - q_i p_i)$, $j \neq i$). So we do not see a ‘freezing’ of the states of the bandits (boxes) that are not searched.

However, consider a different, discounted-reward problem, in which there are no costs, but we receive a reward of $\beta^T/(1 - \beta)$ if the object is found at time T . Note that

$$(1 + \beta + \beta^2 + \dots) - \frac{\beta^T}{1 - \beta} = 1 + \beta + \dots + \beta^{T-1} \rightarrow T, \text{ as } \beta \rightarrow 1.$$

So in the limit $\beta \rightarrow 1$ the problem of maximizing the expected value of the reward $\beta^T/(1 - \beta)$ that we obtain when the object is found at time T , is the same as the problem of minimizing ET .

Now a policy is just a sequence in which we will search the boxes. If the k th search of box i occurs at time t_k this produces expected reward in the new problem of $p_i(1 - q_i)^{k-1}q_i\beta^{t_k}$, i.e. the probability the object is in box i and found on the k th search. So the problem is now the same as that for a SFABP in which the k th time that B_i is continued we get reward of $p_i(1 - q_i)^{k-1}q_i$. This is decreasing in k , which means that the index for B_i when it has been searched $k - 1$ times is simply $p_i(1 - q_i)^{k-1}q_i$ (as we will take $\tau = 1$ in (7.2)). When the boxes has been searched k_1, \dots, k_n times the posterior probabilities are $p'_i \propto p_i(1 - q_i)^{k_i}$. Thus the index of B_i is proportional to $p'_i q_i$, where $p'_i(k_1, \dots, k_n)$ is the posterior probability that the object is in box i .

Thus we have proved that a myopic policy is optimal: always search next the box for which $p'_i(k_1, \dots, k_n)q_i$ is greatest. (Actually, this is a rather complicated way to prove this fact. Instead we prove it using an interchange argument.)

9 Average-cost Programming

The average-cost optimality equation. Policy improvement algorithm.

9.1 Average-cost optimality equation

Suppose that for a stationary Markov policy π , the following limit exists:

$$\lambda(\pi, x) = \lim_{t \rightarrow \infty} \frac{1}{t} E_{\pi} \left[\sum_{\tau=0}^{t-1} c(x_{\tau}, u_{\tau}) \mid x_0 = x \right].$$

We might expect there to be a well-defined notion of an optimal **average cost**, $\lambda(x) = \inf_{\pi} \lambda(\pi, x)$, and that under appropriate assumptions, $\lambda(x) = \lambda$ should not depend on x . Moreover, a reasonable guess is that

$$F_s(x) = s\lambda + \phi(x) + \epsilon(s, x),$$

where $\epsilon(s, x) \rightarrow 0$ as $s \rightarrow \infty$. Here $\phi(x) + \epsilon(s, x)$ reflects a transient due to the initial state. Suppose that the state space and action space are finite. From the optimality equation for the finite horizon problem we have

$$F_s(x) = \min_u \{c(x, u) + E[F_{s-1}(x_1) \mid x_0 = x, u_0 = u]\}. \quad (9.1)$$

So by substituting $F_s(x) \sim s\lambda + \phi(x)$ into (9.1), we obtain

$$s\lambda + \phi(x) \sim \min_u \{c(x, u) + E[(s-1)\lambda + \phi(x_1) \mid x_0 = x, u_0 = u]\}$$

which suggests that the average-cost optimality equation should be:

$$\lambda + \phi(x) = \min_u \{c(x, u) + E[\phi(x_1) \mid x_0 = x, u_0 = u]\}. \quad (9.2)$$

Theorem 9.1. *Suppose there exists a constant λ and bounded function ϕ satisfying (9.2). Let π be the policy which in each state x chooses u to minimize the right hand side. Then λ is the minimal average-cost and π is the optimal stationary policy.*

Proof. Suppose u is chosen by some policy π' . By repeated substitution of (9.2) into itself we have

$$\phi(x) \leq -t\lambda + E_{\pi'} \left[\sum_{\tau=0}^{t-1} c(x_{\tau}, u_{\tau}) \mid x_0 = x \right] + E_{\pi'} [\phi(x_t) \mid x_0 = x].$$

with equality if $\pi' = \pi$. Divide this by t and let $t \rightarrow \infty$. Boundedness of ϕ ensures that $(1/t)E_{\pi'}[\phi(x_t) \mid x_0 = x] \rightarrow 0$. So we obtain

$$0 \leq -\lambda + \lim_{t \rightarrow \infty} \frac{1}{t} E_{\pi'} \left[\sum_{\tau=0}^{t-1} c(x_{\tau}, u_{\tau}) \mid x_0 = x \right],$$

with equality if $\pi' = \pi$. □

So an average-cost optimal policy can be found by looking for a bounded solution to (9.2). Notice that if ϕ is a solution of (9.2) then so is $\phi + (\text{a constant})$, because the (a constant) will cancel from both sides of (9.2). Thus ϕ is undetermined up to an additive constant. In searching for a solution to (9.2) we can therefore pick any state, say \bar{x} , and arbitrarily take $\phi(\bar{x}) = 0$. The function ϕ is called the **relative value function**.

9.2 Example: admission control at a queue

Each day a consultant is presented with the opportunity to take on a new job. The jobs are independently distributed over n possible types and on a given day the offered type is i with probability a_i , $i = 1, \dots, n$. Jobs of type i pay R_i upon completion. Once he has accepted a job he may accept no other job until that job is complete. The probability that a job of type i takes k days is $(1 - p_i)^{k-1} p_i$, $k = 1, 2, \dots$. Which jobs should the consultant accept?

Solution. Let 0 and i denote the states in which he is free to accept a job, and in which he is engaged upon a job of type i , respectively. Then (9.2) is

$$\begin{aligned}\lambda + \phi(0) &= \sum_{i=1}^n a_i \max[\phi(0), \phi(i)], \\ \lambda + \phi(i) &= (1 - p_i)\phi(i) + p_i[R_i + \phi(0)], \quad i = 1, \dots, n.\end{aligned}$$

Taking $\phi(0) = 0$, these have solution $\phi(i) = R_i - \lambda/p_i$, and hence

$$\lambda = \sum_{i=1}^n a_i \max[0, R_i - \lambda/p_i].$$

The left hand side is increasing in λ and the right hand side is decreasing λ . Hence there is a root, say λ^* , and this is the maximal average-reward. The optimal policy takes the form: *accept only jobs for which $p_i R_i \geq \lambda^*$* .

9.3 Value iteration bounds

Value iteration in the average-cost case is based upon the idea that $F_s(x) - F_{s-1}(x)$ approximates the minimal average-cost for large s . For the rest of this lecture we suppose the state space is finite.

Theorem 9.2. *Define*

$$m_s = \min_x \{F_s(x) - F_{s-1}(x)\}, \quad M_s = \max_x \{F_s(x) - F_{s-1}(x)\}. \quad (9.3)$$

Then $m_s \leq \lambda \leq M_s$, where λ is the minimal average-cost.

Proof. Suppose $\pi = f^\infty$ is the average-cost optimal policy over the infinite horizon, taking $u = f(x)$, with average-cost λ . Then

$$\begin{aligned} F_{s-1}(x) + m_s &\leq F_{s-1}(x) + [F_s(x) - F_{s-1}(x)] \\ &= F_s(x) \\ &\leq c(x, f(x)) + E[F_{s-1}(x_1) \mid x_0 = x, u_0 = f(x)]. \end{aligned}$$

We substitute this into itself $t - 1$ times to get

$$F_{s-1}(x) \leq -m_s t + E_\pi \left[\sum_{\tau=0}^{t-1} c(x_\tau, u_\tau) \mid x_0 = x \right] + E_\pi [F_{s-1}(x_t) \mid x_0 = x].$$

Divide by t and let $t \rightarrow \infty$ to get $m_s \leq \lambda$. A bound $\lambda_s \leq M_s$ is found similarly using

$$\begin{aligned} F_{s-1}(x) &\geq -M_s + c(x, u) + E[F_{s-1}(x_1) \mid x_0 = x, u_0 = f_s(x)]. \\ &\geq -M_s t + E_{\pi_s} \left[\sum_{\tau=0}^{t-1} c(x_\tau, u_\tau) \mid x_0 = x \right] + E_{\pi_s} [F_{s-1}(x_t) \mid x_0 = x]. \end{aligned}$$

□

This justifies use of a **value iteration algorithm** in which we calculate F_s until $M_s - m_s \leq \epsilon m_s$. At that point the stationary policy f_s^∞ achieves an average-cost that is within $\epsilon \times 100\%$ of optimal.

9.4 Policy improvement algorithm

In the average-cost case a **policy improvement algorithm** is based on the following observations. Suppose that for a policy $\pi = f^\infty$, we have that λ, ϕ solve

$$\lambda + \phi(x) = c(x, f(x_0)) + E[\phi(x_1) \mid x_0 = x, u_0 = f(x_0)],$$

and there exists a policy $\pi_1 = f_1^\infty$ such that

$$\lambda + \phi(x) \geq c(x, f_1(x_0)) + E[\phi(x_1) \mid x_0 = x, u_0 = f_1(x_0)], \quad (9.4)$$

for all x , and with strict inequality for some x (and thus $f_1 \neq f$). Then following the lines of proof in Theorems 9.1 and 9.2 (repeatedly substituting (9.4) into itself),

$$\lambda \geq \lim_{t \rightarrow \infty} \frac{1}{t} E_{\pi_1} \left[\sum_{\tau=0}^{t-1} c(x_\tau, u_\tau) \mid x_0 = x \right]. \quad (9.5)$$

So π_1 is at least as good as π . If there is no π_1 then π satisfies (9.2) and so π is optimal. This justifies the following **policy improvement algorithm**

(0) Choose an arbitrary stationary policy π_0 . Set $s = 1$.

(1) For stationary policy $\pi_{s-1} = f_{s-1}^\infty$ determine ϕ, λ to solve

$$\lambda + \phi(x) = c(x, f_{s-1}(x)) + E[\phi(x_1) \mid x_0 = x, u_0 = f_{s-1}(x)].$$

This gives a set of linear equations, and so is intrinsically easier to solve than (9.2). The average-cost of π_{s-1} is λ .

(2) Now determine the policy $\pi_s = f_s^\infty$ from

$$f_s(x) = \arg \min_u \{c(x, u) + E[\phi(x_1) \mid x_0 = x, u_0 = u]\},$$

taking $f_s(x) = f_{s-1}(x)$ whenever this is possible. If $\pi_s = \pi_{s-1}$ then we have a solution to (9.2) and so π_{s-1} is optimal. Otherwise π_s is a new policy. By the calculation in (9.5) this has an average-cost no more than λ , so π_s is at least as good as π_{s-1} . We now return to step (1) with $s := s + 1$.

If both the action and state spaces are finite then there are only a finite number of possible stationary policies and so the policy improvement algorithm must find an optimal stationary policy in finitely many iterations. By contrast, the value iteration algorithm only obtains increasingly accurate approximations of λ^* .

Example 9.3. Consider again the example of §9.2. Let us start with a policy π_0 which accept only jobs of type 1. The average-cost of this policy can be found by solving

$$\lambda + \phi(0) = a_1\phi(1) + \sum_{i=2}^n a_i\phi(0),$$

$$\lambda + \phi(i) = (1 - p_i)\phi(i) + p_i[R_i + \phi(0)], \quad i = 1, \dots, n.$$

The solution is $\lambda = a_1 p_1 R_1 / (a_1 + p_1)$, $\phi(0) = 0$, $\phi(1) = p_1 R_1 / (a_1 + p_1)$, and $\phi(i) = R_i - \lambda / p_i$, $i \geq 2$. The first use of step (1) of the policy improvement algorithm will create a new policy π_1 , which improves on π_0 , by accepting jobs for which $\phi(i) = \max\{\phi(0), \phi(i)\}$, i.e. for which $\phi(i) = R_i - \lambda / p_i > 0 = \phi(0)$.

If there are no such i then π_0 is optimal. So we may conclude that π_0 is optimal if and only if $p_i R_i \leq a_1 p_1 R_1 / (a_1 + p_1)$ for all $i \geq 2$.

Policy improvement in the discounted-cost case.

In the case of strict discounting the policy improvement algorithm is similar:

(0) Choose an arbitrary stationary policy π_0 . Set $s = 1$.

(1) For stationary policy $\pi_{s-1} = f_{s-1}^\infty$ determine G to solve

$$G(x) = c(x, f_{s-1}(x)) + \beta E[G(x_1) \mid x_0 = x, u_0 = f_{s-1}(x)].$$

(2) Now determine the policy $\pi_s = f_s^\infty$ from

$$f_s(x) = \arg \min_u \{c(x, u) + \beta E[G(x_1) \mid x_0 = x, u_0 = u]\},$$

taking $f_s(x) = f_{s-1}(x)$ whenever this is possible. Stop if $f_s = f_{s-1}$. Otherwise return to step (1) with $s := s + 1$.

10 Continuous-time Markov Decision Processes

Control problems in a continuous-time stochastic setting. Markov jump processes when the state space is discrete. Uniformization

10.1 Stochastic scheduling on parallel machines

A collection of n jobs is to be processed on a single machine. They have processing times X_1, \dots, X_n , which are *ex ante* distributed as independent exponential random variables, $X_i \sim \mathcal{E}(\lambda_i)$ and $EX_i = 1/\lambda_i$, where $\lambda_1, \dots, \lambda_n$ are known.

If jobs are processed in order $1, 2, \dots, n$, they finished in expected time $1/\lambda_1 + \dots + 1/\lambda_n$. So the order of processing does not matter.

But now suppose there are m (≥ 2) identical machines working in parallel. Let C_i be the **completion time** of job i .

- $\max_i C_i$ is called the **makespan** (the time when all jobs are complete).
- $\sum_i C_i$ is called the **flow time** (sum of completion times).

Suppose we wish to minimize the expected makespan. We can find the optimal order of processing by stochastic dynamic programming. But now we are in continuous time, $t \geq 0$. So we need the important facts:

(i) $\min(X_i, X_j) \sim \mathcal{E}(\lambda_i + \lambda_j)$; (ii) $P(X_i < X_j \mid \min(X_i, X_j) = t) = \lambda_i/(\lambda_i + \lambda_j)$.

Suppose $m = 2$. The optimality equations are

$$\begin{aligned} F(\{i\}) &= \frac{1}{\lambda_i} \\ F(\{i, j\}) &= \frac{1}{\lambda_i + \lambda_j} [1 + \lambda_i F(\{j\}) + \lambda_j F(\{i\})] \\ F(S) &= \min_{i, j \in S} \frac{1}{\lambda_i + \lambda_j} [1 + \lambda_i F(S^i) + \lambda_j F(S^j)], \end{aligned}$$

where S is a set of uncompleted jobs, and we use the abbreviated notation $S^i = S \setminus \{i\}$.

It is helpful to rewrite the optimality equation. Let $\Lambda = \sum_i \lambda_i$. Then

$$\begin{aligned} F(S) &= \min_{i, j \in S} \frac{1}{\Lambda} \left[1 + \lambda_i F(S^i) + \lambda_j F(S^j) + \sum_{k \neq i, j} \lambda_k F(S) \right] \\ &= \min_{\substack{u_i \in [0, 1], i \in S, \\ \sum_i u_i \leq 2}} \frac{1}{\Lambda} \left[1 + \Lambda F(S) + \sum_i u_i \lambda_i (F(S^i) - F(S)) \right] \end{aligned}$$

This is helpful, because in all equations there is now the same divisor, Λ . An event occurs after a time that is exponentially distributed with parameter Λ , but with probability λ_k/Λ this is a ‘dummy event’ if $k \neq i, j$. This trick is known as **uniformization**. Having set this up we might also then say let $\Lambda = 1$.

We see that it is optimal to start by processing the two jobs in S for which $\delta_i(S) := \lambda_i(F(S^i) - F(S))$ is least.

Theorem 10.1.

(a) *Expected makespan is minimized by LHR.*

(b) *Expected flow time is minimized by HHR.*

(c) $E[C_{(n-m+1)}]$ (*expected time there is first an idle machine*) is minimized by LHR.

Proof. We prove only (a), and for simplicity assume $m = 2$ and $\lambda_1 < \dots < \lambda_n$. We would like to prove that for all $i, j \in S \subseteq \{1, \dots, n\}$ that

$$i < j \iff \delta_i(S) < \delta_j(S) \quad (\text{except possibly if both } i \text{ and } j \text{ are the jobs that would be processed by the optimal policy}). \quad (10.1)$$

This would mean that jobs should be started in the order $1, 2, \dots, n$.

Let π be LLR. Take an induction hypothesis that (10.1) is true and that $F(S) = F(\pi, S)$ when S is a strict subset of $\{1, \dots, n\}$. Now consider $S = \{1, \dots, n\}$. We examine $F(\pi, S)$, and $\delta_i(\pi, S)$, under the conjectured optimal policy π . Let S^k denote $S \setminus \{k\}$. For $i \geq 3$,

$$\begin{aligned} F(\pi, S) &= \frac{1}{\lambda_1 + \lambda_2} [1 + \lambda_1 F(S^1) + \lambda_2 F(S^2)] \\ F(\pi, S^i) &= \frac{1}{\lambda_1 + \lambda_2} [1 + \lambda_1 F(S^{1i}) + \lambda_2 F(S^{2i})] \\ \implies \delta_i(\pi, S) &= \frac{1}{\lambda_1 + \lambda_2} [\lambda_1 \delta_i(S^1) + \lambda_2 \delta_i(S^2)], \quad i \geq 3. \end{aligned} \quad (10.2)$$

If for some $3 \leq i < j$ we were to have $\delta_i(\pi, S) > \delta_j(\pi, S)$ then this would require that either $\delta_i(S^1) > \delta_j(S^1)$ or $\delta_i(S^2) > \delta_j(S^2)$. But our inductive hypothesis for (10.1) rules these out.

Similarly, we can compute $\delta_1(\pi, S)$.

$$\begin{aligned} F(\pi, S) &= \frac{1}{\lambda_1 + \lambda_2 + \lambda_3} [1 + \lambda_1 F(S^1) + \lambda_2 F(S^2) + \lambda_3 F(\pi, S)] \\ F(\pi, S^1) &= \frac{1}{\lambda_1 + \lambda_2 + \lambda_3} [1 + \lambda_1 F(S^1) + \lambda_2 F(S^{12}) + \lambda_3 F(S^{13})] \\ \implies \delta_1(\pi, S) &= \frac{1}{\lambda_1 + \lambda_2 + \lambda_3} [\lambda_2 \delta_1(S^2) + \lambda_3 \delta_1(\pi, S) + \lambda_1 \delta_3(S^1)] \\ &= \frac{1}{\lambda_1 + \lambda_2} [\lambda_1 \delta_3(S^1) + \lambda_2 \delta_1(S^2)]. \end{aligned} \quad (10.3)$$

By comparing (10.2) and (10.3) we see that we could only have $\delta_i(S) < \delta_1(S)$ for $i \geq 3$ if at least one of $\delta_i(S^1) < \delta_3(S^1)$ or $\delta_i(S^2) < \delta_1(S^2)$ is true. These are ruled out by our inductive hypothesis. Similarly, we cannot have $\delta_i(S) < \delta_2(S)$ for $i \geq 3$.

This completes a step of a step of an inductive proof by showing that (10.1) is true for S , and that $F(S) = F(\pi, S)$. We only need to check the base of the induction. This is provided by the simple calculation

$$\begin{aligned} \delta_1(\{1, 2\}) &= \lambda_1(F(\{2\}) - F(\{1, 2\})) = \lambda_1 \left[\frac{1}{\lambda_2} - \frac{1}{\lambda_1 + \lambda_2} \left(1 + \frac{\lambda_1}{\lambda_2} + \frac{\lambda_2}{\lambda_1} \right) \right] \\ &= -\frac{\lambda_2}{\lambda_1 + \lambda_2} \leq \delta_2(\{1, 2\}). \end{aligned} \quad \square$$

The proof of (b) is very similar, except that the inequality in (10.1) should be reversed. The base of the induction comes from $\delta_1(\{1, 2\}) = -1$.

The proof of (c) is also similar. The base of the induction is provided by $\delta_1(\{1, 2\}) = \lambda_1(0 - 1/(\lambda_1 + \lambda_2))$. Since we are seeking to maximize $EC_{(n-m+1)}$ we should process jobs for which δ_i is greatest, i.e., least λ_i . The problem in (c) is known as the **Lady's nylon stocking problem**. We think of a lady (having $m = 2$ legs) who starts with n stockings, wears two at a time, each of which may fail, and she wishes to maximize the expected time until she has only one good stocking left to wear.

10.2 Controlled Markov jump processes

The above is an example of a controlled **Markov jump process**, which evolves in continuous time within a discrete state space. In general:

- The state is i . We choose some control, say u .
- After a time that is exponentially distributed with parameter $q_i(u) = \sum_{j \neq i} q_{ij}(u)$, (i.e. having mean $1/q_i(u)$), the state jumps.
- Until the jump occurs cost accrues at rate $c(i, u)$.
- The jump is to state j ($\neq i$) with probability $q_{ij}(u)/q_i(u)$.

The infinite-horizon optimality equation is

$$F(i) = \min_u \left\{ \frac{1}{q_i(u)} \left[c(i, u) + \sum_j q_{ij}(u) F(j) \right] \right\}.$$

We can use the **uniformization** trick, picking a B such that $q_i(u) \leq B$ for all i, u .

$$F(i) = \min_u \left\{ \frac{1}{B} \left[c(i, u) + (B - q_i(u)) F(i) + \sum_j q_{ij}(u) F(j) \right] \right\}.$$

We now have something that looks exactly like a discrete-time optimality equation

$$F(i) = \min_u \left\{ \bar{c}(i, u) + \sum_j p_{ij}(u) F(j) \right\}$$

where $\bar{c}(i, u) = c(i, u)/B$, $p_{ij}(u) = q_{ij}(u)/B$, $j \neq i$, and $P_{ii}(u) = 1 - q_i(u)/B$.

This is great! It means we can use all the methods and theorems that we have developed previously for solving discrete-time dynamic programming problems.

We can also introduce discounting by imagining that there is an ‘exponential clock’ of rate α which takes the state to a place where no further cost or reward is obtained. This leads to an optimality equation of the form

$$F(i) = \min_u \left\{ \bar{c}(i, u) + \beta \sum_j p_{ij}(u) F(j) \right\},$$

where $\beta = B/(B + \alpha)$, $\bar{c}(i, u) = c(i, u)/(B + \alpha)$, and $p_{ij}(u)$ is as above.

10.3 Example: admission control at a queue

Consider a queue in which the number of customers waiting may be $0, 1, \dots$. There is a constant service rate μ (meaning that the service times of customers are distributed $\mathcal{E}(\mu)$, and arrival rate u , where u is controllable between 0 and a maximum value M). Let $c(x, u) = ax - Ru$. This corresponds to paying a **holding cost** a per unit time for each customer in the queue and receiving a reward R at the point that each new customer is admitted (and therefore incurring reward at rate Ru when the arrival rate is u). Suppose there is discounting at rate α . The problem is one of choosing $0 \leq u_t \leq M$ to minimize

$$E \left[\int_0^\infty (ax_t - Ru_t) e^{-\alpha t} dt \right].$$

Let us take $B = \alpha + M + \mu$, and without loss of generality assume $B = 1$.

After uniformization the discounted-cost optimality equation will look like

$$F(0) = \inf_{u \in [0, M]} \{-Ru + \beta[uF(1) + (M - u)F(j - 1)]\}$$

$$F(x) = \inf_{u \in [0, M]} \{ax - Ru + \beta[uF(x + 1) + \mu F(x - 1) + (M - u)F(x)]\}, \quad x = 1, 2, \dots$$

So we can see that the optimal control is bang-bang, taking $u = 0$ or $u = M$ as the coefficient of u , namely $-R + F(x + 1) - F(x)$, is positive or negative. One can set up a value iteration form of this, i.e.

$$F_{k+1}(0) = \inf_{u \in [0, M]} \{-Ru + \beta[uF_k(1) + (M - u)F_k(j - 1)]\}$$

$$F_{k+1}(x) = \inf_{u \in [0, M]} \{ax - Ru + \beta[uF_k(x + 1) + \mu F_k(x - 1) + (M - u)F_k(x)]\}, \quad x = 1, 2, \dots$$

and then prove by induction that $F_k(x)$ is concave in x . This means that there exists a **threshold rule** such that the optimal policy will be of the form:

$$u = \begin{cases} 0 & \geq x^* \\ M & < x^* \end{cases}.$$

Time-average cost optimality. The optimality equation is

$$\begin{aligned}\phi(0) + \gamma &= \inf_{u \in [0, M]} [-Ru + u\phi(1) + (\mu + M - u)\phi(0)], \\ &= \inf_{u \in [0, M]} [u\{-R + \phi(1) - \phi(0)\} + (\mu + M)\phi(0)],\end{aligned}$$

$$\begin{aligned}\phi(x) + \gamma &= \inf_{u \in [0, M]} [ax - Ru + u\phi(x+1) + \mu\phi(x-1) + (M - u)\phi(x)], \\ &= \inf_{u \in [0, M]} [ax + u\{-R + \phi(x+1) - \phi(x)\} + \mu\phi(x-1) + M\phi(x)], \quad x > 0.\end{aligned}$$

Thus u should be chosen to be 0 or 1 as $-R + \phi(x+1) - \phi(x)$ is positive or negative.

Let us consider what happens under the policy that take $u = M$ for all x . The relative costs for this policy, say f , are given by

$$f(x) + \gamma = ax - R\lambda + Mf(x+1) + \mu f(x-1), \quad x > 0.$$

The solution to the homogeneous part of this recursion is of the form $f(x) = d_1 1^x + d_2 (\mu/M)^x$. Assuming $M < \mu$ and we desire a solution for f that does not grow exponentially, we take $d_2 = 0$ and so the solution is effectively the solution to the inhomogeneous part, i.e.

$$f(x) = \frac{ax(x+1)}{2(\mu - M)}, \quad \gamma = \frac{aM}{\mu - M} - MR,$$

Applying the idea of policy improvement, we conclude that a better policy is to take $u = 0$ (i.e. don't admit a customer) if $-R + f(x+1) - f(x) > 0$, i.e. if

$$\frac{(x+1)a}{\mu - M} - R > 0.$$

Further policy improvement would probably be needed to reach the optimal policy. However, this policy already exhibits an interesting property: it rejects customers for smaller queue length x than does a policy which rejects a customer if and only if

$$\frac{(x+1)a}{\mu} - R > 0.$$

This second policy is optimal if one is purely concerned with whether or not an individual customer that joins when there are x customers in front of him will show a profit on the basis of the difference between the reward R and his expected holding cost $(x+1)a/\mu$. This example exhibits the difference between **individual optimality** (which is myopic) and **social optimality**. The socially optimal policy is more reluctant to admit customers because it anticipates that more customers are on the way; thus it feels less badly about forgoing the profit on a customer that presents himself now, recognizing that admitting such a customer can cause customers who are admitted after him to suffer greater delay. As expected, the policies are nearly the same if the arrival rate λ is small.

11 Restless Bandits

11.1 Examples

Again, we start with a family of n alternative Markov decision processes. Given their states at time t , say $x_1(t), \dots, x_n(t)$, we are to choose actions $u_1(t), \dots, u_n(t)$ to apply at time t . As in the multi-armed bandit problem, we suppose that there are just two available actions, so $u_i(t) \in \{0, 1\}$. We generalize the SFABP set up in two ways.

Our first generalization is to require that at each time t exactly m ($< n$) of the bandits be given the ‘active’ action $u_i = 1$.

Our second generalization is that the ‘passive’ action $u = 0$ no longer freezes a bandit; instead, the state evolves, but differently from its continuation under $u = 1$.

Example 11.1. Suppose the state of a bandit is a measure of its vigour. The active and passive actions correspond to notions of work and rest. Performance is positively related to vigour (or lack of fatigue), which increases with rest and decreases with work. For example, suppose that x takes values in $\{1, \dots, k\}$. If the active action $u = 1$ is applied to a bandit in state x , then there accrues an immediate reward of $r(x)$, increasing in x , but vigour decreases to $\max\{x - 1, 1\}$. The passive action $u = 0$ produces no reward, but vigour increases to $\min\{x + 1, k\}$.

Example 11.2. The active and passive actions might correspond to notions of ‘observation’ and ‘no observation’. Suppose that each bandit is in one of two conditions: 0 and 1, associated with being ‘bad’ or ‘good’, respectively. It moves back and forth between these two conditions independently of any actions applied by the decision-maker, according to a 2-state Markov chain. Each bandit is now a POMDP. So far as the decision-maker is concerned the state of the bandit is the probability that it is in good condition. Under the action $u = 1$ the condition is observed, and if this is found to be i then $x(t + 1) = p_{i1}$. Moreover, if the condition is good then a reward is obtained. Under the action $u = 0$ the underlying condition of the process is not observed, and so, in a Bayesian manner, $x(t + 1) = x(t)p_{11} + (1 - x(t))p_{01}$. No reward is obtained.

Example 11.3. The active and passive actions correspond to running the process at different speeds. For example, suppose for $0 < \epsilon < 1$,

$$P(j|i, 0) = \begin{cases} \epsilon P(j|i, 1), & i \neq j \\ (1 - \epsilon) + \epsilon P(i|i, 1), & i = j \end{cases}$$

Thus a bandit which is operated continuously with $u = 1$ has the same stationary distribution as one that is operated continuously with $u = 0$. But the process moves faster when $u = 1$.

11.2 Whittle index policy

Let $\Omega = \{(u_1, \dots, u_n) : u_i \in \{0, 1\} \text{ for all } i, \text{ and } \sum_i u_i(t) = m\}$. The optimality equation is

$$F(x) = \max_{u \in \Omega} \left\{ \sum_i r(x_i, u_i) + \beta \sum_{y_1, \dots, y_n} F(y_1, \dots, y_n) \prod_i P(y_i | x_i, u_i) \right\}.$$

Let us focus upon average reward (i.e. the limit as $\beta \rightarrow 1$). This is attractive because performance does not depend on the initial state. Assuming that the n bandits are statistically equivalent it is plausible that, under an optimal policy, bandit i will be given the action $u_i = 1$ for precisely a fraction m/n of the time. This motivates interest in an upper bound on the maximal average reward that can be obtained by considering a single bandit and asking how it should be controlled if we wish to maximize the average reward obtained from that bandit, subject to a relaxed constraint that $u_i = 1$ is employed for a fraction of exactly m/n of the time.

So consider a stationary Markov policy for operating a single restless bandit. Let z_x^u be the proportion of time that the bandit is in state x and that under this policy the action u is taken. An upper bound for our problem can be found from a linear program in variables $\{z_x^u : x \in E, u \in \{0, 1\}\}$:

$$\text{maximize } \sum_{x,u} r(x, u) z_x^u \tag{11.1}$$

subject to

$$\sum_{x,u} z_x^u = 1 \tag{11.2}$$

$$\sum_x z_x^0 \geq 1 - m/n \tag{11.3}$$

$$\sum_u z_x^u = \sum_{y,u} z_y^u P(x | y, u), \text{ for all } x \tag{11.4}$$

$$z_x^u \geq 0, \text{ for all } x, u. \tag{11.5}$$

Here (11.4) are equations that determine the stationary probabilities. Notice that we have put an inequality in (11.3). Let us justify this by making the assumption that action $u = 1$ (which we call the **active action**) is in some sense better than $u = 0$ (which we call the **passive action**). So if constraint (11.3) did not exist then we would wish to take $u = 1$ in all states. At optimality (11.3) will hold with equality.

The optimal value of the **dual LP** problem is equal to g , where this can be found from the average reward dynamic programming equation

$$\phi(x) + g = \max_{u \in \{0,1\}} \left\{ r(x, u) + \lambda(1 - u) + \sum_y \phi(y) P(y | x, u(x)) \right\}. \tag{11.6}$$

Here λ and $\phi(x)$ are the Lagrange multipliers for constraints (11.3) and (11.4), respectively. The multiplier λ is positive and may be interpreted as a *subsidy* for taking the passive action. It is interesting to see how (11.6) can be obtained from (11.1)–(11.4). However, we might have simply taken as our starting point a problem of maximizing average reward when there is a subsidy for taking the passive action.

In general, the solution of (11.6) partitions the state space E into three sets, E_0 , E_1 and E_{01} , where, respectively the optimal action is $u = 0$, $u = 1$, or some randomization between both $u = 0$ and $u = 1$. Let us avoid uninteresting pathologies by supposing that the state space is finite, and that every pure policy gives rise to a Markov chain with one recurrent class. Then the set E_{01} , where there is randomization, need never contain more than 1 state, a fact that is known for general Markov decision processes with constraints.

It is reasonable to expect that as the subsidy λ increases in (11.6) the set of states E_0 (in which $u = 0$ is optimal) should increase monotonically. This need not happen in general. However, if it does then we say the bandit is **indexable**. Whittle defines as an index the least value of the subsidy λ such that $u = 0$ is optimal. We call this the **Whittle index**, denoting it $W(\cdot)$, where $W(x) = \inf\{\lambda : x \in E_0(\lambda)\}$. It can be used to define a heuristic policy (the **Whittle index policy**) in which, at each instant, one engages m bandits with the greatest indices, i.e. those that are the last to leave the set E_1 as the subsidy for the passive action increases. The Whittle index extends the Gittins optimal policy for classical bandits; Whittle indices can be computed separately for each bandit; they are the same as the Gittins index in the case that $u = 0$ is a freezing action, so that $P(j|i, 0) = \delta_{ij}$.

11.3 Whittle indexability

The discussion so far begs two questions: (i) under what assumptions is a restless bandit indexable, and (ii) how good is the Whittle index policy? Might it be optimal, or very nearly optimal as n becomes large?

Interestingly, there are special classes of restless bandit for which one can prove indexability. Bandits of the type in Example 11.2 are indexable. The dual-speed bandits in Example 11.3 are indexable. A restless bandit is also indexable if the passive action transition probabilities, $P(j | i, 0)$, depend only on j (the destination state).

11.4 Fluid models of large stochastic systems

It is often interesting to think about problems in some ‘large N ’ limit. Consider, for example, N identical independently running single server queues, of type $M/M/1$, each with its own Poisson arrival stream of rate λ and server of rate μ . The probability that a given queue has i customers is $\pi_i = \rho^i(1 - \rho)$ where $\rho = \lambda/\mu$. The queues are running independently and so we would expect the number of them that have i customers to be $N\pi_i$. Suppose we start off with $Nx_i(0)$ of the queues having i customers, where $\sum_i x_i(0) = 1$. Since N is large, transitions will be happening very fast, and so using

the law of large numbers we expect to see

$$\begin{aligned}\frac{d}{dt}x_0(t) &= \mu x_1(t) - \lambda x_0(t) \\ \frac{d}{dt}x_i(t) &= \lambda x_{i-1}(t) + \mu x_{i+1}(t) - (\lambda + \mu)x_i(t).\end{aligned}$$

We have replaced our stochastic system by a deterministic fluid approximation. (There are theorems which talk about the convergence when $N \rightarrow \infty$.) These differential equations will produce a trajectory $x_i(t) \rightarrow \pi_i$ as $t \rightarrow \infty$.

The same thing happens even if we link the behaviour of the queues. Suppose we have total processing effort $N\mu$. Rather than place μ per queue, as above, we now decide to allocate 2μ to all queues having more than the median number of customers. Suppose $\sum_{i=1}^{j-1} x_i(t) + \alpha x_j(t) = 0.5$, so queues with $\leq j$ customers are not served, those with $\geq j+1$ are served, and some service $(1-\alpha)x_j(t)\mu$ effort is wasted. Now the fluid approximation is

$$\begin{aligned}\frac{d}{dt}x_0(t) &= \mu x_1(t) - \lambda x_0(t) \\ \frac{d}{dt}x_i(t) &= \lambda x_{i-1}(t) - \lambda x_i(t), \quad 0 < i \leq j-1 \\ \frac{d}{dt}x_j(t) &= \lambda x_{j-1}(t) + 2\mu x_{j+1}(t) - \lambda x_j(t) \\ \frac{d}{dt}x_k(t) &= \lambda x_{k-1}(t) + 2\mu x_{k+1}(t) - (\lambda + 2\mu)x_k(t), \quad k \geq j+1.\end{aligned}$$

The appropriate set of differential equations will of course change depending upon which j is the median queue size. There will still be convergence to some equilibrium point (which we might hope will have a smaller average queue size.)

11.5 Asymptotic optimality

We now turn to the question of optimality or near optimality of the Whittle index policy. Taking $m = \alpha n$, let $R_W^{(n)}(\alpha)$, $R_{\text{opt}}^{(n)}(\alpha)$ and $r(\alpha)$ denote, respectively, the average reward that is obtained from n restless bandits under the Whittle index policy, under an optimal policy, and from a single bandit under the relaxed policy (that the bandit receive the action $u = 1$ for a fraction α of the time). Then

$$R_W^{(n)}(\alpha) \leq R_{\text{opt}}^{(n)}(\alpha) \leq nr(\alpha).$$

It is plausible that the Whittle index policy should be **asymptotically optimal** as n becomes large, in the sense that $r(\alpha) - R_W^{(n)}(\alpha)/n \rightarrow 0$ as $n \rightarrow \infty$. This is true if certain differential equations have an asymptotically stable equilibrium point (i.e. a point to which they converge from any starting state). These are the differential equations which describe a fluid approximation to the stochastic process of the bandit states evolving under the Whittle index policy.

Suppose bandits move on a state space of size k and let $z_i(t)$ be the proportion of the bandits in state i . The ‘fluid approximation’ for large n is given by piecewise linear differential equations, of the form:

$$\frac{dz}{dt} = A(z)x + b(z),$$

where $A(z)$ and $b(z)$ are constants within k polyhedral regions which partition the positive orthant of \mathbb{R}^k . For example for $k = 2$,

$$\frac{dz_i}{dt} = \sum_j q_{ji}(z)z_j - \sum_j q_{ij}(z)z_i$$

$$\frac{dz_1}{dt} = \begin{cases} -(q_{12}^0 + q_{21}^0)z_1 + (q_{12}^0 - q_{12}^1)\rho + q_{21}^0, & z_1 \geq \rho \\ -(q_{12}^1 + q_{21}^1)z_1 - (q_{21}^0 - q_{21}^1)\rho + q_{21}^0, & z_1 \leq \rho \end{cases}$$

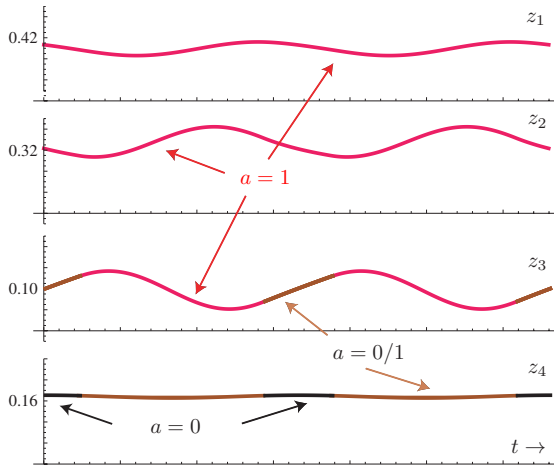
It turns out that there are examples of dual-speed restless bandits (needing $k > 3$) in which the differential equations have an asymptotically stable equilibrium cycle (rather than a stable equilibrium point), and this can lead to suboptimality of the Whittle index policy. However, in examples, the suboptimality was never found to be more than about one part in 10^4 .

Theorem 11.4. *If bandits are indexable, and the fluid model has an asymptotically stable equilibrium point, then the Whittle index heuristic is asymptotically optimal, — in the sense that the reward per bandit tends to the reward that is obtained under the relaxed policy.*

Here is an example where the Whittle index does not quite provide asymptotically optimal performance.

$$(q_{ij}^0) = \begin{pmatrix} -2 & 1 & 0 & 1 \\ 2 & -2 & 0 & 0 \\ 0 & 56 & -\frac{113}{2} & \frac{1}{2} \\ 1 & 1 & \frac{1}{2} & -\frac{5}{2} \end{pmatrix}, \quad (q_{ij}^1) = \begin{pmatrix} -2 & 1 & 0 & 1 \\ 2 & -2 & 0 & 0 \\ 0 & \frac{7}{25} & -\frac{113}{400} & \frac{1}{400} \\ 1 & 1 & \frac{1}{2} & -\frac{5}{2} \end{pmatrix}$$

$$r^0 = (0, 1, 10, 10), \quad r^1 = (10, 10, 10, 0), \quad \rho = 0.835$$



Equilibrium point is $(\bar{z}_1, \bar{z}_2, \bar{z}_3, \bar{z}_4) = (0.409, 0.327, 0.100, 0.164)$. $\bar{z}_1 + \bar{z}_2 + \bar{z}_3 = 0.836$. The equilibrium is a cycle. Relaxed policy obtains average reward 10 per bandit. Heuristic obtains only 9.9993 per bandit.

It is tempting to try to generalize the idea of a Whittle index for restless bandits to problems with a discounted reward criterion, starting with the appropriate functional equation in place of and adding a subsidy for use of the passive action. However, there is no asymptotic optimality result for this case that is analogous to the result of for the average reward case. The use of discounted versions of Whittle indices can actually end up recommending the worst of all priority policies, and a payoff that is very far from the optimum. This is because the identity of the optimal priority policy can critically depend on the starting state of the n restless bandits, whereas the ordering of Whittle indices is calculated without reference to the starting state.

12 Sequential Assignment and Allocation Problems

Having met the Secretary problem, Multi-armed bandit problem, etc., we now turn to some other very well-known and interesting problems that I have personally enjoyed.

12.1 Sequential stochastic assignment problem

Derman, Lieberman and Ross (1974) defined the following **sequential stochastic assignment problem** (SSAP). It has been applied in many contexts, including kidney transplantation, aviation security, buying decisions in supply chains, and real estate.

There are n workers available to perform n jobs. First job 1 appears, followed by job 2, etc. Associated with the j th job is a random variable X_j which takes the value x_j . X_1, \dots, X_n are i.i.d. random variables with distribution function G . If a ‘perfect’ worker is assigned to the value x_j job, a reward x_j is obtained. However, none of the workers is perfect, and whenever the i th worker is assigned to any type x_j job, the (expected) reward is given by $p_i x_j$, where $0 < p_i < 1$ is a known constant. After a worker is assigned, he is unavailable for future assignments. The problem is to assign the n workers to the n jobs so as to maximize the total expected reward. A policy is a rule for assigning workers to jobs. Let random variable i_j be the worker (identified by number) assigned to the j th arriving job. The expected reward to be maximized is

$$\sum_{j=1}^n E[p_{i_j} X_j].$$

The optimal policy is given by the following theorem. The surprise is that the thresholds $\alpha_{i,n}$ are independent of the p 's.

Theorem 12.1. *For each $n > 1$, there exist numbers $a_{0,n} > a_{1,n} > a_{2,n} > \dots > a_{n,n} = 0$ such that whenever there are n stages to go and $p_1 > \dots > p_n$ then the optimal choice in the initial stage is to use p_i if the random variable X_1 is contained in the interval $(a_{i-1,n}, a_{i,n}]$. The $a_{i,n}$ depend on G but are independent of the p_i .*

Furthermore $a_{i,n}$ is the expected value, in an $n - 1$ stage problem, of the quantity to which the i th largest p is assigned (assuming an optimal policy is followed), and

$$F(p_1, \dots, p_{n-1}) = \sum_{i=1}^{n-1} p_i a_{i,n}, \quad p_1 > \dots > p_{n-1}.$$

Proof. The proof is by induction. Assuming it is true for $n - 1$, we have for n ,

$$\begin{aligned} F(p_1, \dots, p_n \mid x) &= \max_i \{x p_k + F(p_1, \dots, p_{k-1}, p_{k+1}, \dots, p_n)\} \\ &= \max_k \left\{ \sum_{j=1}^{k-1} p_j a_{j,n} + p_k x + \sum_{j=k+1}^n p_j a_{j-1,n} \right\}. \end{aligned}$$

We use the **Hardy-Littlewood rearrangement inequality**, which says that if $a < A$ and $b < B$ then $ab + AB > aB + bA$. That is, it is best to match smallest with smallest, largest with largest.

Suppose $X_1 = x_1$ and $\{a_{1,n}, \dots, a_{i-1,n}, x_1, a_{i,n}, \dots, a_{n-1,n}\}$ is a decreasing sequence. Then the optimum matching of these numbers against the decreasing numbers $\{p_1, p_2, \dots, p_n\}$, is to form the sum of the products obtained by matching, for each j , the j th largest of $\{a_{1,n}, \dots, a_{i-1,n}, x_1, a_{i,n}, \dots, a_{n-1,n}\}$ with j th largest of $\{p_1, p_2, \dots, p_n\}$, which means that x_1 should be matched with p_i . Notice that $a_{k,n+1}$ is the expected value of the X that gets matched to p_k . This value does not depend on the values of the p_i s. \square

We could implement the optimal strategy by offering each job to the workers in decreasing order of their p values. Worker i will accept the job if workers $1, \dots, i - 1$ reject it, and then $X \geq a_{i,n}$, since the job is worth as much to him as he would expect to get if he forgoes it and then faces a $n - 1$ stage problem (where his expected match is $a_{i,n}$). This is nice. We can obtain the socially optimal allocation by presenting the workers with a problem that they each solve from an individually optimal viewpoint.

12.2 Sequential allocation problems

Groundwater Management Burt (1965). Each day water is to be pumped from an aquifer, and replenished by a random amount of rainfall R_t . The aim is to maximize an expected sum of utilities $\sum a(y_t)$ minus pumping cost $\sum c(x_t, y_t)$.

$$F(x, s) = \max_{y \in [0, x]} \{a(y) - c(x, y) + \beta EF(x - y + R_s, s - 1)\}$$

$$F(x, 0) = 0. \quad x \text{ is level of water in an aquifer.}$$

Here $s = h - t$ is time-to-go.

Investment problem Derman, Lieberman and Ross (1975). With probability p_t there is an opportunities to invest part of ones capital at time t . The aim is maximize the expected sum of $\sum a(y_t)$.

$$F(x, s) = q_s F(x, s - 1) + p_s \max_{y \in [0, x]} \{a(y) + F(x - y, s - 1)\}$$

$$F(x, 0) = 0. \quad x \text{ is remaining capital of dollars.}$$

General fighter problem With probability p_t there is an opportunity for a fighter to shoot down an enemy plane. If m missiles are used then the enemy plane is destroyed with probability $a(m)$ and the fighter survives the dogfight with probability $c(m)$. The aim is to maximize the expected number of enemy planes destroyed.

$$F(n, s) = q_s F(n, s - 1) + p_s \max_{m \in \{1, \dots, n\}} \{a(m) + c(m)F(n - m, s - 1)\}$$

$$F(n, 0) = 0. \quad n \text{ is remaining stock of missiles.}$$

Bomber problem Klinger and Brown (1968). With probability p_t a bomber must defend itself against an attack and wishes to maximize the probability of reaching its final target.

$$\begin{aligned} P(n, s) &= P(\text{survive to for } s \text{ further distance}) \\ &= q_s P(n, s - 1) + p_s \max_{m \in \{1, \dots, n\}} c(m) P(n - m, s - 1) \end{aligned}$$

$$P(n, 0) = 1. \quad n \text{ is remaining stock of } \text{missiles}. \text{ Typically, } c(m) = 1 - \theta^m.$$

The bomber problem can also be posed in continuous time, as

$$P(n, t) = e^{-t} + \int_0^t \max_{m \in \{1, \dots, n\}} c(m) P(n - m, s) e^{-(t-s)} ds.$$

Intuitively obvious properties of an optimal policy are (for the bomber problem, and other problems similarly)

- (A) $m(n, s) \searrow$ as $s \nearrow$
- (B) $m(n, s) \nearrow$ as $n \nearrow$
- (C) $n - m(n, s) \nearrow$ as $n \nearrow$

Properties like these are sometimes quite easy to prove. Sometimes this is by a value iteration approaching, proving that the value function has appropriate concavity properties. Or sometimes an interchange arguments helps. Consider (C) for the Bomber Problem. We shall assume that $\log c(m)$ is concave in m .

Proof of (C) of the Bomber Problem. Let

$$m(n, s) = \arg \max_{m \in \{1, \dots, n\}} c(m) P(n - m, s - 1).$$

We wish to show that $n - m(n, s - 1)$ is nondecreasing in n . Suppose this were not the case. So perhaps $m = m(n, s - 1)$ but $m' = m(n + 1, s - 1)$ with $n - m > n + 1 - m'$, i.e., $m' > m + 1$. Consider the product of the survival probabilities

$$c(m) P(n - m, s - 1) \times c(m') P(n + 1 - m', s - 1) \tag{12.1}$$

Let $\bar{m} = m' - 1$ and $\bar{m}' = m + 1$. This different choice of amounts to fire in state $(n, s - 1)$ and $(n + 1, s - 1)$ would have a product of survival probabilities

$$\begin{aligned} &c(\bar{m}) P(n - \bar{m}, s - 1) \times c(\bar{m}') P(n + 1 - \bar{m}', s - 1) \\ &= c(m' - 1) P(n + 1 - m', s - 1) \times c(m + 1) P(n - m, s - 1) \end{aligned} \tag{12.2}$$

$$(12.2) - (12.1)$$

$$\begin{aligned} &= \left[c(m + 1) c(m' - 1) - c(m) c(m') \right] P(n - m, s - 1) P(n + 1 - m', s - 1) \\ &\geq 0, \end{aligned}$$

since $\log c(m)$ is concave means that $m' > m + 1 \implies \frac{c(m+1)}{c(m)} > \frac{c(m')}{c(m'-1)}$.

Hence at least one of our original m and m' must not have been optimal. \square

However **(B)** for the Bomber problem remains an unproven conjecture, as also is **(A)** for the General fighter problem. It has been shown there are **(B)** is not true for the Bomber problem if $c(m)$ is an arbitrary concave function. However, opinion is very divided about whether or not **(B)** might be true for the special concave function $c(m) = 1 - \theta^m$. Very extensive computation have turned up no counterexample.

12.3 SSAP with arrivals

Suppose workers and jobs arrive according to Poisson processes with respective rates γ and λ . The workers are identical with $p_i = 1$ (the so-called house-selling case). Rewards are exponentially discounted with rate α . Job values are i.i.d., say $U[0, 1]$.

Suppose that an arriving job is offered to the workers in inverse order of their arrival times; so the worker that arrived most recently has first right of refusal for jobs, and workers try to maximize their own expected job values. The individually optimal (IO) policy is the Nash equilibrium of a noncooperative game; that is, if all workers follow the IO policy and each worker is trying to maximize its own expected job value, then no worker will have incentive to deviate from the IO policy. Righter (2011) has shown that the IO policy is unique and has proved the following.

Theorem 12.2. *The IO policy is socially optimal (maximizing total discounted return). Thus the socially optimal policy is a threshold policy.*

A worker who is i th to be offered a job of value x should accept it iff $x \geq t_i$, where $t_1 > t_2 > \dots$ and t_i is the expected discounted job value that is allocated to worker i under the IO policy. Use uniformization, so $\gamma + \lambda + \alpha = 1$. This mean we can now think of γ and λ as probabilities, and of our problem as taking place in discrete time. The thresholds are given by thinking about worker i being indifferent between accepting and rejecting a job:

$$t_i = \underbrace{\gamma t_{i+1}}_{\text{new worker arrives}} + \lambda \left[\underbrace{t_i P(X < t_i)}_{\text{new job assigned to worker behind } i} + \underbrace{E[X 1_{\{t_i \leq X < t_{i-1}\}}]}_{\text{new job assigned to worker } i} + \underbrace{t_{i-1} P(X \geq t_{i-1})}_{\text{new job assigned to worker ahead of } i} \right].$$

Proof. We show that following the IO policy, starting with the first decision, is better (for the sum of the worker's obtained values) than following an arbitrary policy for the first decision and then switching to the IO policy thereafter (call the latter policy π). Essentially we are showing that the policy improvement algorithm cannot improve IO.

Suppose at the first decision, $t_{i-1} > x \geq t_i$ but the job is assigned to no worker by π , and the IO policy is used thereafter. Workers $1, 2, \dots, i-1$ and all future arriving

workers will have the same expected job value (EJV) under IO and π . Worker i would have had x under IO, but will only get t_i under π . Workers $i + 1, \dots, n$ will also have greater EJVs from time 1 onward once i has been assigned a job, so for them also IO is better than π . So the job should be assigned.

What if there are n workers present and $x < t_n$? IO rejects the job. If a policy π assigns the job to a worker, which we may take to be worker n , then all other workers have the same EVJ under IO and π , but worker n is taking x , whereas under π his EJV would be t_n , which is greater. So π cannot be optimal. \square

12.4 SSAP with a postponement option

Consider now a SSAP with m perfect workers ($p_i = 1$), and $n (> m)$ jobs to be presented sequentially, with i.i.d. values $X_1, \dots, X_n \sim U[0, 1]$, and discounting. We no longer demand that a job must be assigned or rejected upon first seeing its value. It may be held in a queue, for possible assignment later. For a state in which there are $s = n - t$ jobs still to see, m workers still unassigned, and a queue holding jobs of values $x_1 > \dots > x_m$ (some of these can be 0) the optimality equations are

$$F_{s,m}(x_1, \dots, x_m) = \int_0^1 G_{s,m}(T(x, x_1, \dots, x_m)) dx$$

$$G_{s,m}(x_1, \dots, x_m) = \max_{i \in \{0, \dots, m\}} \left\{ \sum_{j=1}^i x_j + \beta F_{s,m-i}(x_{i+1}, \dots, x_m) \right\},$$

where $T(x, x_1, \dots, x_m)$ is the vector formed from $\{x, x_1, \dots, x_m\}$ by discarding the smallest element and then rearranging the rest in decreasing order. We have written the optimality equation in two parts. The first equation is about receiving the next job. The second is about assigning jobs that are presently in the postponement queue. Feng and Hartman (2012) have proved the following.

Theorem 12.3. *An optimal policy is to assign the available job of greatest value, x , iff $x \geq \alpha_{m,s}$ where the threshold $\alpha_{m,s}$ depends on the numbers of unassigned workers, m , and jobs left to see, s , but not on the values of the other jobs in the queue.*

The final part of this statement is rather surprising. One might have expected that it would be optimal to assign the jobs of greatest value if and only if it is greater than some threshold, but where this is a complicated function of m , s and the values of the other jobs that are in the queue and are available to be assigned later.

The proof is very complicated. It would be nice to find a simpler proof, perhaps by thinking about re-casting the problem into one about individual optimality, as Righter did with the arrivals case of the SSAP.

12.5 Stochastic knapsack and bin packing problems

Similar to the problems we have addressed thus far is the stochastic **knapsack problem**. It has been applied within the field of revenue management. The capacity of the knapsack is a given amount of resource that can be used to fulfill customer demand over a certain time frame. Some examples include: rooms in a hotel aimed at weekend tourists, or seats on an airplane that must be sold before departure. Items arrive with random sizes X_i and values R_i . We wish to maximize the expected total value of the items we can fit into a knapsack of size c .

For example, Coffman, Flatto, Weber (1987) have studied a stochastic **bin packing problem**. Items of sizes X_1, X_2, \dots, X_n are encountered sequentially. These are i.i.d. with distribution function G . We wish to maximize the expected number of items that we can pack into a bin of size c . Each item must be accepted or rejected at the time it is encountered.

Let x be the space that is left in the bin. The dynamic programming equation is

$$F_s(x) = (1 - G(x))F_{s-1}(x) + \int_0^x \max\{F_{s-1}(x), 1 + F_{s-1}(x - y)\} dG(y)$$
$$F_0(x) = 0.$$

The optimal rule is to accept the sth to last item iff doing so would leave remaining space in the bin of at least $z_{s-1,x}$, where $F_{s-1}(z_{s-1,x}) + 1 = F_{s-1}(x)$.

13 LQ Regulation

Models with linear dynamics and quadratic costs in discrete and continuous time. Riccati equation, and its validity with additive white noise. Linearization of nonlinear models.

13.1 The LQ regulation problem

As we have seen, the elements of a control optimization problem are specification of (i) the dynamics of the process, (ii) which quantities are observable at a given time, and (iii) an optimization criterion.

In the **LQG model** the plant equation and observation relations are **linear**, the cost is **quadratic**, and the noise is **Gaussian** (jointly normal). The LQG model is important because it has a complete theory and illuminates key concepts, such as controllability, observability and the certainty-equivalence principle.

Begin with a model in which the state x_t is fully observable and there is no noise. The plant equation of the time-homogeneous $[A, B, \cdot]$ system has the linear form

$$x_t = Ax_{t-1} + Bu_{t-1}, \quad (13.1)$$

where $x_t \in \mathbb{R}^n$, $u_t \in \mathbb{R}^m$, A is $n \times n$ and B is $n \times m$. The cost function is

$$\mathbf{C} = \sum_{t=0}^{h-1} c(x_t, u_t) + \mathbf{C}_h(x_h), \quad (13.2)$$

with one-step and terminal costs

$$c(x, u) = x^\top R x + u^\top S x + x^\top S^\top u + u^\top Q u = \begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} R & S^\top \\ S & Q \end{pmatrix} \begin{pmatrix} x \\ u \end{pmatrix}, \quad (13.3)$$

$$\mathbf{C}_h(x) = x^\top \Pi_h x. \quad (13.4)$$

All quadratic forms are non-negative definite ($\succeq 0$), and Q is positive definite ($\succ 0$). There is no loss of generality in assuming that R , Q and Π_h are symmetric. This is a model for **regulation** of (x, u) to the point $(0, 0)$ (i.e. steering to a critical value).

To solve the optimality equation we shall need the following lemma.

Lemma 13.1. *Suppose x, u are vectors. Consider a quadratic form*

$$\begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} \Pi_{xx} & \Pi_{xu} \\ \Pi_{ux} & \Pi_{uu} \end{pmatrix} \begin{pmatrix} x \\ u \end{pmatrix}.$$

Assume it is symmetric and $\Pi_{uu} > 0$, i.e. positive definite. Then the minimum with respect to u is achieved at

$$u = -\Pi_{uu}^{-1} \Pi_{ux} x,$$

and is equal to

$$x^\top [\Pi_{xx} - \Pi_{xu} \Pi_{uu}^{-1} \Pi_{ux}] x.$$

Proof. Suppose the quadratic form is minimized at u . Then

$$\begin{aligned} & \begin{pmatrix} x \\ u+h \end{pmatrix}^\top \begin{pmatrix} \Pi_{xx} & \Pi_{xu} \\ \Pi_{ux} & \Pi_{uu} \end{pmatrix} \begin{pmatrix} x \\ u+h \end{pmatrix} \\ &= x^\top \Pi_{xx} x + 2x^\top \Pi_{xu} u + \underbrace{2h^\top \Pi_{ux} x + 2h^\top \Pi_{uu} u}_{\text{linear term}} + u^\top \Pi_{uu} u + h^\top \Pi_{uu} h. \end{aligned}$$

To be stationary at u , the underbraced linear term in h^\top must be zero, so

$$u = -\Pi_{uu}^{-1} \Pi_{ux} x,$$

and the optimal value is $x^\top [\Pi_{xx} - \Pi_{xu} \Pi_{uu}^{-1} \Pi_{ux}] x$. \square

Theorem 13.2. *Assume the structure of (13.1)–(13.4). Then the value function has the quadratic form*

$$F(x, t) = x^\top \Pi_t x, \quad t \leq h, \quad (13.5)$$

and the optimal control has the linear form

$$u_t = K_t x_t, \quad t < h.$$

The time-dependent matrix Π_t satisfies the Riccati equation

$$\Pi_t = f \Pi_{t+1}, \quad t < h, \quad (13.6)$$

where Π_h has the value given in (13.4), and f is an operator having the action

$$f \Pi = R + A^\top \Pi A - (S^\top + A^\top \Pi B)(Q + B^\top \Pi B)^{-1} (S + B^\top \Pi A). \quad (13.7)$$

The $m \times n$ matrix K_t is given by

$$K_t = -(Q + B^\top \Pi_{t+1} B)^{-1} (S + B^\top \Pi_{t+1} A), \quad t < h. \quad (13.8)$$

Proof. Assertion (13.5) is true at time h . Assume it is true at time $t+1$. Then

$$\begin{aligned} F(x, t) &= \inf_u [c(x, u) + (Ax + Bu)^\top \Pi_{t+1} (Ax + Bu)] \\ &= \inf_u \left[\begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} R + A^\top \Pi_{t+1} A & S^\top + A^\top \Pi_{t+1} B \\ S + B^\top \Pi_{t+1} A & Q + B^\top \Pi_{t+1} B \end{pmatrix} \begin{pmatrix} x \\ u \end{pmatrix} \right]. \end{aligned}$$

Lemma 13.1 shows the minimizer is $u = K_t x$, and gives the form of f . \square

13.2 The Riccati recursion

The backward recursion (13.6)–(13.7) is called the **Riccati equation**.

(i) Since the optimal control is linear in the state, say $u = Kx$, an equivalent expression for the Riccati equation is

$$f\Pi = \inf_K [R + K^\top S + S^\top K + K^\top QK + (A + BK)^\top \Pi(A + BK)].$$

(ii) The optimally controlled process obeys $x_{t+1} = \Gamma_t x_t$. We call Γ_t the **gain matrix** and it is given by

$$\Gamma_t = A + BK_t = A - B(Q + B^\top \Pi_{t+1} B)^{-1} (S + B^\top \Pi_{t+1} A).$$

(iii) S can be normalized to zero by choosing a new control $u^* = u + Q^{-1}Sx$, and setting $A^* = A - BQ^{-1}S$, $R^* = R - S^\top Q^{-1}S$. So $A^*x + Bu^* = Ax + Bu$ and $c(x, u) = x^\top Rx + u^{*\top} Q u^*$.

(iv) Similar results are true if $x_{t+1} = A_t x_t + B_t u_t + \alpha_t$, where $\{\alpha_t\}$ is a known sequence of disturbances, and the aim is to track a sequence of values (\bar{x}_t, \bar{u}_t) , $t = 0, \dots, h-1$, so the cost is

$$c(x, u, t) = \begin{pmatrix} x - \bar{x}_t \\ u - \bar{u}_t \end{pmatrix}^\top \begin{pmatrix} R_t & S_t^\top \\ S_t & Q_t \end{pmatrix} \begin{pmatrix} x - \bar{x}_t \\ u - \bar{u}_t \end{pmatrix}.$$

13.3 White noise disturbances

Suppose the plant equation (13.1) is now

$$x_{t+1} = Ax_t + Bu_t + \epsilon_t,$$

where $\epsilon_t \in \mathbb{R}^n$ is vector **white noise**, defined by the properties $E\epsilon = 0$, $E\epsilon_t \epsilon_t^\top = N$ and $E\epsilon_t \epsilon_s^\top = 0$, $t \neq s$. The dynamic programming equation is then

$$F(x, t) = \inf_u \{c(x, u) + E_\epsilon [F(Ax + Bu + \epsilon, t + 1)]\},$$

with $F(x, h) = x^\top \Pi_h x$. Try a solution $F(x, t) = x^\top \Pi_t x + \gamma_t$. This holds for $t = h$. Suppose it is true for $t + 1$, then

$$\begin{aligned} F(x, t) &= \inf_u \{c(x, u) + E(Ax + Bu + \epsilon)^\top \Pi_{t+1} (Ax + Bu + \epsilon) + \gamma_{t+1}\} \\ &= \inf_u \{c(x, u) + (Ax + Bu)^\top \Pi_{t+1} (Ax + Bu) \\ &\quad + 2E\epsilon^\top \Pi_{t+1} (Ax + Bu)\} + E[\epsilon^\top \Pi_{t+1} \epsilon] + \gamma_{t+1} \\ &= \inf_u \{c(x, u) + (Ax + Bu)^\top \Pi_{t+1} (Ax + Bu)\} + \text{tr}(N\Pi_{t+1}) + \gamma_{t+1}, \end{aligned}$$

where $\text{tr}(A)$ means the trace of matrix A . Here we use the fact that

$$E[\epsilon^\top \Pi \epsilon] = E\left[\sum_{ij} \epsilon_i \Pi_{ij} \epsilon_j\right] = E\left[\sum_{ij} \epsilon_j \epsilon_i \Pi_{ij}\right] = \sum_{ij} N_{ji} \Pi_{ij} = \text{tr}(N\Pi).$$

Thus (i) Π_t follows the same Riccati equation as before, (ii) the optimal control is $u_t = K_t x_t$, and (iii)

$$F(x, t) = x^\top \Pi_t x + \gamma_t = x^\top \Pi_t x + \sum_{j=t+1}^h \text{tr}(N \Pi_j).$$

The final term can be viewed as the cost of correcting future noise. In the infinite horizon limit of $\Pi_t \rightarrow \Pi$ as $t \rightarrow \infty$, we incur an average cost per unit time of $\text{tr}(N \Pi)$, and a transient cost of $x^\top \Pi x$ that is due to correcting the initial x .

13.4 LQ regulation in continuous-time

In continuous-time we take $\dot{x} = Ax + Bu$ and cost

$$\mathbf{C} = \int_0^h \begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} R & S^\top \\ S & Q \end{pmatrix} \begin{pmatrix} x \\ u \end{pmatrix} dt + (x^\top \Pi x)_h.$$

We can obtain the continuous-time solution from the discrete time solution by moving forward in time in increments of Δ . Make the following replacements.

$$x_{t+1} \rightarrow x_{t+\Delta}, \quad A \rightarrow I + A\Delta, \quad B \rightarrow B\Delta, \quad R, S, Q \rightarrow R\Delta, S\Delta, Q\Delta.$$

Then as before, $F(x, t) = x^\top \Pi x$, where Π obeys the Riccati equation

$$\frac{\partial \Pi}{\partial t} + R + A^\top \Pi + \Pi A - (S^\top + \Pi B)Q^{-1}(S + B^\top \Pi) = 0.$$

This is slightly simpler than the discrete time version. The optimal control is $u(t) = K(t)x(t)$, where $K(t) = -Q^{-1}(S + B^\top \Pi)$.

The optimally controlled plant equation is $\dot{x} = \Gamma(t)x$, where

$$\Gamma(t) = A + BK = A - BQ^{-1}(S + B^\top \Pi).$$

13.5 Linearization of nonlinear models

Linear models are important because they arise naturally via the linearization of nonlinear models. Consider the state-structured nonlinear model:

$$\dot{x} = a(x, u).$$

Suppose x, u are perturbed from an equilibrium (\bar{x}, \bar{u}) where $a(\bar{x}, \bar{u}) = 0$. Let $x' = x - \bar{x}$ and $u' = u - \bar{u}$. The linearized version is

$$\dot{x}' = \dot{x} = a(\bar{x} + x', \bar{u} + u') = Ax' + Bu'$$

where

$$A_{ij} = \left. \frac{\partial a_i}{\partial x_j} \right|_{(\bar{x}, \bar{u})}, \quad B_{ij} = \left. \frac{\partial a_i}{\partial u_j} \right|_{(\bar{x}, \bar{u})}.$$

If (\bar{x}, \bar{u}) is to be a stable equilibrium point then we must be able to choose a control that can bring the system back to (\bar{x}, \bar{u}) from any nearby starting point.

14 Controllability and Observability

Controllability in discrete and continuous time. Stabilizability.

14.1 Controllability and Observability

The discrete-time system $[A, B, \cdot]$ is defined by the plant equation

$$x_t = Ax_{t-1} + Bu_{t-1}, \quad (14.1)$$

The **controllability** question is: can we bring x to an arbitrary prescribed value by some u -sequence?

The discrete-time system $[A, B, C]$ is defined by (14.1) and observation relation

$$y_t = Cx_{t-1}. \quad (14.2)$$

$y_t \in \mathbb{R}^p$ is observed, but x_t is not. C is $p \times n$. The **observability** question is: can we infer x_0 from subsequent y values?

Definition 14.1. The $[A, B, \cdot]$ system is **r-controllable** if one can bring it from an arbitrary prescribed x_0 to an arbitrary prescribed x_r by some u -sequence u_0, u_1, \dots, u_{r-1} . A system of dimension n is **controllable** if it is r -controllable for some r

Definition 14.2. The $[A, B, C]$ system is said to be **r-observable** if x_0 can be inferred from knowledge of the observations y_1, \dots, y_r and relevant control values u_0, \dots, u_{r-2} for any initial x_0 . An n -dimensional system is **observable** if r -observable for some r .

The notion of observability stands in dual relation to that of controllability; a duality that indeed persists throughout the subject.

14.2 Controllability

Example 14.3. Consider the case, ($n = 2, m = 1$),

$$x_t = \begin{pmatrix} a_{11} & 0 \\ a_{21} & a_{22} \end{pmatrix} x_{t-1} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} u_{t-1}.$$

This system is not 1-controllable. But

$$x_2 - A^2 x_0 = Bu_1 + ABu_0 = \begin{pmatrix} 1 & a_{11} \\ 0 & a_{21} \end{pmatrix} \begin{pmatrix} u_1 \\ u_0 \end{pmatrix}.$$

So it is 2-controllable if and only if $a_{21} \neq 0$.

In general, by substituting the plant equation (14.1) into itself, we see that we must find u_0, u_1, \dots, u_{r-1} to satisfy

$$\Delta = x_r - A^r x_0 = Bu_{r-1} + ABu_{r-2} + \dots + A^{r-1}Bu_0, \quad (14.3)$$

for arbitrary Δ . In providing conditions for controllability we use the following theorem.

Theorem 14.4. (The Cayley-Hamilton theorem) Any $n \times n$ matrix A satisfies its own characteristic equation. So that if

$$\det(\lambda I - A) = \sum_{j=0}^n a_j \lambda^{n-j}$$

then $\sum_{j=0}^n a_j A^{n-j} = 0$.

The implication is that $I, A, A^2, \dots, A^{n-1}$ contains basis for A^r , $r = 0, 1, \dots$. We are now in a position to characterize controllability.

Theorem 14.5. (i) The system $[A, B, \cdot]$ is r -controllable iff the matrix

$$M_r = [B \quad AB \quad A^2B \quad \dots \quad A^{r-1}B]$$

has rank n , (ii) equivalently, iff the $n \times n$ matrix

$$M_r M_r^\top = \sum_{j=0}^{r-1} A^j (B B^\top) (A^\top)^j$$

is nonsingular (or, equivalently, positive definite.) (iii) If the system is r -controllable then it is s -controllable for $s \geq \min(n, r)$, and (iv) a control transferring x_0 to x_r with minimal cost $\sum_{t=0}^{r-1} u_t^\top u_t$ is

$$u_t = B^\top (A^\top)^{r-t-1} (M_r M_r^\top)^{-1} (x_r - A^r x_0), \quad t = 0, \dots, r-1.$$

Proof. (i) The system (14.3) has a solution for arbitrary Δ iff M_r has rank n .

(ii) That is, iff there does not exist nonzero w such that $w^\top M_r = 0$. Now

$$\exists w \neq 0 : w^\top M_r = 0 \iff \exists w \neq 0 : w^\top M_r M_r^\top w = 0 \iff M_r M_r^\top \text{ is not p.d.s.}$$

(iii) The rank of M_r is non-decreasing in r , so if the system is r -controllable, it is $(r+1)$ -controllable. By the Cayley-Hamilton theorem, the rank is constant for $r \geq n$.

(iv) Consider the Lagrangian

$$\sum_{t=0}^{r-1} u_t^\top u_t + \lambda^\top \left(\Delta - \sum_{t=0}^{r-1} A^{r-t-1} B u_t \right),$$

giving

$$u_t = \frac{1}{2} B^\top (A^\top)^{r-t-1} \lambda.$$

Now we can determine λ from (14.3). □

14.3 Controllability in continuous-time

Theorem 14.6. (i) The n dimensional system $[A, B, \cdot]$ is controllable iff the matrix M_n has rank n , or (ii) equivalently, iff

$$G(t) = \int_0^t e^{As} B B^\top e^{A^\top s} ds,$$

is positive definite for all $t > 0$. (iii) If the system is controllable then a control that achieves the transfer from $x(0)$ to $x(t)$ with minimal control cost $\int_0^t u_s^\top u_s ds$ is

$$u(s) = B^\top e^{A^\top(t-s)} G(t)^{-1} (x(t) - e^{At} x(0)).$$

Note that there is now no notion of r -controllability. However, $G(t) \downarrow 0$ as $t \downarrow 0$, so the transfer becomes more difficult and costly as $t \downarrow 0$.

14.4 Example: broom balancing

Consider the problem of balancing a broom in an upright position on your hand. By Newton's laws, the system obeys $m(\ddot{u} \cos \theta + L\ddot{\theta}) = mg \sin \theta$.

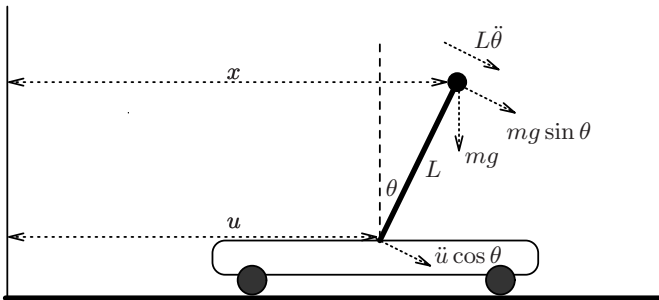


Figure 1: Force diagram for broom balancing

For small θ we have $\cos \theta \sim 1$ and $\theta \sim \sin \theta = (x - u)/L$. So with $\alpha = g/L$

$$\ddot{x} = \alpha(x - u),$$

equivalently,

$$\frac{d}{dt} \begin{pmatrix} x \\ \dot{x} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ \alpha & 0 \end{pmatrix} \begin{pmatrix} x \\ \dot{x} \end{pmatrix} + \begin{pmatrix} 0 \\ -\alpha \end{pmatrix} u.$$

Since

$$[B \quad AB] = \begin{bmatrix} 0 & -\alpha \\ -\alpha & 0 \end{bmatrix},$$

the system is controllable if θ is initially small.

14.5 Stabilizability

Suppose we apply the stationary closed-loop control $u = Kx$ so that $\dot{x} = Ax + Bu = (A + BK)x$. So with $\Gamma = A + BK$, we have

$$\dot{x} = \Gamma x, \quad x_t = e^{\Gamma t} x_0, \quad \text{where } e^{\Gamma t} = \sum_{j=0}^{\infty} (\Gamma t)^j / j!$$

Similarly, in discrete-time, we have can take the stationary control, $u_t = Kx_t$, so that $x_t = Ax_{t-1} + Bu_{t-1} = (A + BK)x_{t-1}$. Now $x_t = \Gamma^t x_0$.

We are interested in choosing Γ so that $x_t \rightarrow 0$ and $t \rightarrow \infty$.

Definition 14.7.

Γ is a **stability matrix** in the continuous-time sense if all its eigenvalues have negative real part, and hence $x_t \rightarrow 0$ as $t \rightarrow \infty$.

Γ is a **stability matrix** in the discrete-time sense if all its eigenvalues of lie strictly inside the unit disc in the complex plane, $|z| = 1$, and hence $x_t \rightarrow 0$ as $t \rightarrow \infty$.

The $[A, B]$ system is said to **stabilizable** if there exists a K such that $A + BK$ is a stability matrix.

Note that $u_t = Kx_t$ is linear and Markov. In seeking controls such that $x_t \rightarrow 0$ it is sufficient to consider only controls of this type since, as we see in the next lecture, such controls arise as optimal controls for the infinite-horizon LQ regulation problem.

14.6 Example: pendulum

Consider a pendulum of length L , unit mass bob and angle θ to the vertical. Suppose we wish to stabilise θ to zero by application of a force u . Then

$$\ddot{\theta} = -(g/L) \sin \theta + u.$$

We change the state variable to $x = (\theta, \dot{\theta})$ and write

$$\frac{d}{dt} \begin{pmatrix} \theta \\ \dot{\theta} \end{pmatrix} = \begin{pmatrix} \dot{\theta} \\ -(g/L) \sin \theta + u \end{pmatrix} \sim \begin{pmatrix} 0 & 1 \\ -g/L & 0 \end{pmatrix} \begin{pmatrix} \theta \\ \dot{\theta} \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} u.$$

Suppose we try to stabilise with a control that is a linear function of only θ (not $\dot{\theta}$), so $u = Kx = (-\kappa, 0)x = -\kappa\theta$. Then

$$\Gamma = A + BK = \begin{pmatrix} 0 & 1 \\ -g/L & 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} (-\kappa \quad 0) = \begin{pmatrix} 0 & 1 \\ -g/L - \kappa & 0 \end{pmatrix}.$$

The eigenvalues of Γ are $\pm \sqrt{-g/L - \kappa}$. So either $-g/L - \kappa > 0$ and one eigenvalue has a positive real part, in which case there is in fact instability, or $-g/L - \kappa < 0$ and eigenvalues are purely imaginary, which means we will in general have oscillations. So successful stabilization must be a function of $\dot{\theta}$ as well, (and this would come out of solution to the LQ regulation problem.)

14.7 Example: satellite in a plane orbit

Consider a satellite of unit mass in a planar orbit and take polar coordinates (r, θ) .

$$\ddot{r} = r\dot{\theta}^2 - \frac{c}{r^2} + u_r, \quad \ddot{\theta} = -\frac{2\dot{r}\dot{\theta}}{r} + \frac{1}{r}u_\theta,$$

where u_r and u_θ are the radial and tangential components of thrust. If $u_r = u_\theta = 0$ then there is a possible equilibrium in which the orbit is a circle of radius $r = \rho$, $\dot{\theta} = \omega = \sqrt{c/\rho^3}$ and $\dot{r} = \dot{\theta} = 0$.

Consider a perturbation of this orbit and measure the deviations from the orbit by

$$x_1 = r - \rho, \quad x_2 = \dot{r}, \quad x_3 = \theta - \omega t, \quad x_4 = \dot{\theta} - \omega.$$

Then, with $n = 4$, $m = 2$,

$$\dot{x} \sim \begin{pmatrix} 0 & 1 & 0 & 0 \\ 3\omega^2 & 0 & 0 & 2\omega\rho \\ 0 & 0 & 0 & 1 \\ 0 & -2\omega/\rho & 0 & 0 \end{pmatrix} x + \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 1/\rho \end{pmatrix} \begin{pmatrix} u_r \\ u_\theta \end{pmatrix} = Ax + Bu.$$

It is easy to check that $M_2 = [B \ AB]$ has rank 4 and that therefore the system is controllable.

Suppose $u_r = 0$ (radial thrust fails). Then

$$B = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1/\rho \end{bmatrix} \quad M_4 = [B \ AB \ A^2B \ A^3B] = \begin{bmatrix} 0 & 0 & 2\omega & 0 \\ 0 & 2\omega & 0 & -2\omega^3 \\ 0 & 1/\rho & 0 & -4\omega^2/\rho \\ 1/\rho & 0 & -4\omega^2/\rho & 0 \end{bmatrix}.$$

which is of rank 4, so the system is still controllable. We can change the radius by tangential braking or thrust.

But if $u_\theta = 0$ (tangential thrust fails). Then

$$B = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad M_4 = [B \ AB \ A^2B \ A^3B] = \begin{bmatrix} 0 & 1 & 0 & -\omega^2 \\ 1 & 0 & -\omega^2 & 0 \\ 0 & 0 & -2\omega/\rho & 0 \\ 0 & -2\omega/\rho & 0 & 2\omega^3/\rho \end{bmatrix}.$$

Since $(2\omega\rho, 0, 0, \rho^2)M_4 = 0$, this is singular and has only rank 3. In fact, the uncontrollable component is the angular momentum, $2\omega\rho\delta r + \rho^2\delta\dot{\theta} = \delta(r^2\dot{\theta})|_{r=\rho, \dot{\theta}=\omega}$.

15 Observability and the LQG Model

LQ regulation problem over the infinite horizon. More on observability. Least squares estimation and the LQG model.

15.1 Infinite horizon limits

Consider the time-homogeneous case and write the finite-horizon cost in terms of time to go s . The terminal cost, when $s = 0$, is denoted $F_0(x) = x^\top \Pi_0 x$. In all that follows we take $S = 0$, without loss of generality.

Lemma 15.1. *Suppose $\Pi_0 = 0$, $R \succeq 0$, $Q \succeq 0$ and $[A, B, \cdot]$ is controllable or stabilizable. Then $\{\Pi_s\}$ has a finite limit Π .*

Proof. Costs are non-negative, so $F_s(x)$ is non-decreasing in s . Now $F_s(x) = x^\top \Pi_s x$. Thus $x^\top \Pi_s x$ is non-decreasing in s for every x . To show that $x^\top \Pi_s x$ is bounded we use one of two arguments.

If the system is controllable then $x^\top \Pi_s x$ is bounded because there is a policy which, for any $x_0 = x$, will bring the state to zero in at most n steps and at finite cost and can then hold it at zero with zero cost thereafter.

If the system is stabilizable then there is a K such that $\Gamma = A + BK$ is a stability matrix. Using $u_t = Kx_t$, we have $x_t = \Gamma^t x$ and $u_t = K\Gamma^t x$, so

$$F_s(x) \leq \sum_{t=0}^{\infty} (x_t^\top R x_t + u_t^\top Q u_t) = x^\top \left[\sum_{t=0}^{\infty} (\Gamma^\top)^t (R + K^\top Q K) \Gamma^t \right] x < \infty.$$

Hence in either case we have an upper bound and so $x^\top \Pi_s x$ tends to a limit for every x . By considering $x = e_j$, the vector with a unit in the j th place and zeros elsewhere, we conclude that the j th element on the diagonal of Π_s converges. Then taking $x = e_j + e_k$ it follows that the off diagonal elements of Π_s also converge. \square

Both value iteration and policy improvement are effective ways to compute the solution to an infinite-horizon LQ regulation problem. Policy improvement goes along the lines developed in Lecture 9.

15.2 Observability

From (14.1) and (14.2) we can determine y_t in terms of x_0 and subsequent controls:

$$x_t = A^t x_0 + \sum_{s=0}^{t-1} A^s B u_{t-s-1},$$
$$y_t = C x_{t-1} = C \left[A^{t-1} x_0 + \sum_{s=0}^{t-2} A^s B u_{t-s-2} \right].$$

Thus, if we define the ‘reduced observation’

$$\tilde{y}_t = y_t - C \left[\sum_{s=0}^{t-2} A^s B u_{t-s-2} \right],$$

then x_0 is to be determined from the system of equations

$$\tilde{y}_t = CA^{t-1}x_0, \quad 1 \leq t \leq r. \quad (15.1)$$

By hypothesis, these equations are mutually consistent, and so have a solution; the question is whether this solution is unique. This is the reverse of the situation for controllability, when the question was whether the equation for u had a solution at all, unique or not. Note that an implication of the system definition is that the property of observability depends only on the matrices A and C ; not upon B at all.

Theorem 15.2. (i) *The system $[A, \cdot, C]$ is r -observable iff the matrix*

$$N_r = \begin{bmatrix} C \\ CA \\ CA^2 \\ \vdots \\ CA^{r-1} \end{bmatrix}$$

has rank n , or (ii) equivalently, iff the $n \times n$ matrix

$$N_r^\top N_r = \sum_{j=0}^{r-1} (A^\top)^j C^\top C A^j$$

is nonsingular. (iii) If the system is r -observable then it is s -observable for $s \geq \min(n, r)$, and (iv) the determination of x_0 can be expressed

$$x_0 = (N_r^\top N_r)^{-1} \sum_{j=1}^r (A^\top)^{j-1} C^\top \tilde{y}_j. \quad (15.2)$$

Proof. If the system has a solution for x_0 (which is so by hypothesis) then this solution must be unique iff the matrix N_r has rank n , whence assertion (i). Assertion (iii) follows from (i). The equivalence of conditions (i) and (ii) can be verified directly as in the case of controllability.

If we define the deviation $\eta_t = \tilde{y}_t - CA^{t-1}x_0$ then the equation amounts to $\eta_t = 0$, $1 \leq t \leq r$. If these equations were not consistent we could still define a ‘least-squares’ solution to them by minimizing any positive-definite quadratic form in these deviations with respect to x_0 . In particular, we could minimize $\sum_{t=0}^{r-1} \eta_t^\top \eta_t$. This minimization gives (15.2). If equations (15.1) indeed have a solution (i.e. are mutually consistent, as we suppose) and this is unique then expression (15.2) must equal this solution; the actual value of x_0 . The criterion for uniqueness of the least-squares solution is that $N_r^\top N_r$ should be nonsingular, which is also condition (ii). \square

We have again found it helpful to bring in an optimization criterion in proving (iv); this time, not so much to construct one definite solution out of many, but to construct a ‘best-fit’ solution where an exact solution might not have existed. This approach lies close to the statistical approach necessary when observations are corrupted by noise.

15.3 Observability in continuous-time

Theorem 15.3. (i) *The n -dimensional continuous-time system $[A, \cdot, C]$ is observable iff the matrix N_n has rank n , or (ii) equivalently, iff*

$$H(t) = \int_0^t e^{A^\top s} C^\top C e^{As} ds$$

is positive definite for all $t > 0$. (iii) If the system is observable then the determination of $x(0)$ can be written

$$x(0) = H(t)^{-1} \int_0^t e^{A^\top s} C^\top \tilde{y}(s) ds,$$

where

$$\tilde{y}(t) = y(t) - \int_0^t C e^{A(t-s)} B u(s) ds.$$

15.4 Example: satellite in planar orbit

Recall the linearised equation $\dot{x} = Ax$, for perturbations of the orbit of a satellite, (here taking $\rho = 1$), where

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} r - \rho \\ \dot{r} \\ \theta - \omega t \\ \dot{\theta} - \omega \end{pmatrix} \quad A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 3\omega^2 & 0 & 0 & 2\omega \\ 0 & 0 & 0 & 1 \\ 0 & -2\omega & 0 & 0 \end{pmatrix}.$$

By taking $C = [0 \ 0 \ 1 \ 0]$ we see that the system is observable on the basis of angle measurements alone, but not observable for $\tilde{C} = [1 \ 0 \ 0 \ 0]$, i.e. on the basis of radius movements alone.

$$N_4 = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & -2\omega & 0 & 0 \\ -6\omega^3 & 0 & 0 & -4\omega^2 \end{bmatrix} \quad \tilde{N}_4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 3\omega^2 & 0 & 0 & 2\omega \\ 0 & -\omega^2 & 0 & 0 \end{bmatrix}$$

15.5 Imperfect state observation with noise

The full LQG model, whose description has been deferred until now, assumes linear dynamics, quadratic costs and Gaussian noise. Imperfect observation is the most im-

portant point. The model is

$$x_t = Ax_{t-1} + Bu_{t-1} + \epsilon_t, \quad (15.3)$$

$$y_t = Cx_{t-1} + \eta_t, \quad (15.4)$$

where ϵ_t is process noise. The state observations are degraded in that we observe only the p -vector $y_t = Cx_{t-1} + \eta_t$, where η_t is observation noise. Typically $p < n$. In this $[A, B, C]$ system A is $n \times n$, B is $n \times m$, and C is $p \times n$. Assume

$$\text{cov} \begin{pmatrix} \epsilon \\ \eta \end{pmatrix} = E \begin{pmatrix} \epsilon \\ \eta \end{pmatrix} \begin{pmatrix} \epsilon \\ \eta \end{pmatrix}^\top = \begin{pmatrix} N & L \\ L^\top & M \end{pmatrix}$$

and that $x_0 \sim N(\hat{x}_0, V_0)$. Let $W_t = (Y_t, U_{t-1}) = (y_1, \dots, y_t; u_0, \dots, u_{t-1})$ denote the observed history up to time t . Of course we assume that $t, A, B, C, N, L, M, \hat{x}_0$ and V_0 are also known; W_t denotes what might be different if the process were rerun.

Lemma 15.4. *Suppose x and y are jointly normal with zero means and covariance matrix*

$$\text{cov} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} V_{xx} & V_{xy} \\ V_{yx} & V_{yy} \end{bmatrix}.$$

Then the distribution of x conditional on y is Gaussian, with

$$E(x | y) = V_{xy}V_{yy}^{-1}y, \quad (15.5)$$

and

$$\text{cov}(x | y) = V_{xx} - V_{xy}V_{yy}^{-1}V_{yx}. \quad (15.6)$$

Proof. Both y and $x - V_{xy}V_{yy}^{-1}y$ are linear functions of x and y and therefore they are Gaussian. From $E[(x - V_{xy}V_{yy}^{-1}y)y^\top] = 0$ it follows that they are uncorrelated and this implies they are independent. Hence the distribution of $x - V_{xy}V_{yy}^{-1}y$ conditional on y is identical with its unconditional distribution, and this is Gaussian with zero mean and the covariance matrix given by (15.6) \square

The estimate of x in terms of y defined as $\hat{x} = Hy = V_{xy}V_{yy}^{-1}y$ is known as the **linear least squares estimate** of x in terms of y . Even without the assumption that x and y are jointly normal, this linear function of y has a smaller covariance matrix than any other unbiased estimate for x that is a linear function of y . In the Gaussian case, it is also the maximum likelihood estimator.

16 Kalman Filter and Certainty Equivalence

The Kalman filter. Certainty equivalence. Separation principle. Hamilton-Jacobi-Bellman equation. Harvesting fish.

16.1 The Kalman filter

Notice that both x_t and y_t can be written as a linear functions of the unknown noise and the known values of u_0, \dots, u_{t-1} .

$$\begin{aligned} x_t &= A^t x_0 + A^{t-1} B u_0 + \dots + B u_{t-1} + A^{t-1} \epsilon_0 + \dots + A \epsilon_{t-1} + \epsilon_t \\ y_t &= C \left(A^{t-1} x_0 + A^{t-2} B u_0 + \dots + B u_{t-2} + A^{t-2} \epsilon_0 + \dots + A \epsilon_{t-2} + \epsilon_{t-1} \right) + \eta_t \end{aligned}$$

Thus the distribution of x_t conditional on $W_t = (Y_t, U_{t-1})$ must be normal, with some mean \hat{x}_t and covariance matrix V_t . Notice that V_t is policy independent (does not depend on u_0, \dots, u_{t-1}).

The following theorem describes recursive updating relations for \hat{x}_t and V_t .

Theorem 16.1. (The Kalman filter) *Suppose that conditional on W_0 , the initial state x_0 is distributed $N(\hat{x}_0, V_0)$ and the state and observations obey the recursions of the LQG model (15.3)–(15.4). Then conditional on W_t , the current state is distributed $N(\hat{x}_t, V_t)$. The conditional mean and variance obey the updating recursions*

$$\hat{x}_t = A \hat{x}_{t-1} + B u_{t-1} + H_t (y_t - C \hat{x}_{t-1}), \quad (16.1)$$

where the time-dependent matrix V_t satisfies a Riccati equation

$$V_t = g V_{t-1}, \quad t < h,$$

where V_0 is given, and g is the operator having the action

$$gV = N + A V A^\top - (L + A V C^\top)(M + C V C^\top)^{-1}(L^\top + C V A^\top). \quad (16.2)$$

The $p \times m$ matrix H_t is given by

$$H_t = (L + A V_{t-1} C^\top)(M + C V_{t-1} C^\top)^{-1}. \quad (16.3)$$

Compare this to the very similar statement of Theorem 13.2. Notice that (16.2) computes V_t forward in time ($V_t = g V_{t-1}$), whereas (13.7) computes Π_t backward in time ($\Pi_t = f \Pi_{t+1}$).

Proof. The proof is by induction on t . Consider the moment when u_{t-1} has been determined but y_t has not yet observed. The distribution of (x_t, y_t) conditional on (W_{t-1}, u_{t-1}) is jointly normal with means

$$\begin{aligned} E(x_t \mid W_{t-1}, u_{t-1}) &= A \hat{x}_{t-1} + B u_{t-1}, \\ E(y_t \mid W_{t-1}, u_{t-1}) &= C \hat{x}_{t-1}. \end{aligned}$$

Let $\Delta_{t-1} = \hat{x}_{t-1} - x_{t-1}$, which by an inductive hypothesis is $N(0, V_{t-1})$. Consider the **innovations**

$$\begin{aligned}\xi_t &= x_t - E(x_t | W_{t-1}, u_{t-1}) = x_t - (A\hat{x}_{t-1} + Bu_{t-1}) = \epsilon_t - A\Delta_{t-1}, \\ \zeta_t &= y_t - E(y_t | W_{t-1}, u_{t-1}) = y_t - C\hat{x}_{t-1} = \eta_t - C\Delta_{t-1}.\end{aligned}$$

Conditional on (W_{t-1}, u_{t-1}) , these quantities are normally distributed with zero means and covariance matrix

$$\text{cov} \begin{bmatrix} \epsilon_t - A\Delta_{t-1} \\ \eta_t - C\Delta_{t-1} \end{bmatrix} = \begin{bmatrix} N + AV_{t-1}A^\top & L + AV_{t-1}C^\top \\ L^\top + CV_{t-1}A^\top & M + CV_{t-1}C^\top \end{bmatrix} = \begin{bmatrix} V_{\xi\xi} & V_{\xi\zeta} \\ V_{\zeta\xi} & V_{\zeta\zeta} \end{bmatrix}.$$

Thus it follows from Lemma 15.4 that the distribution of ξ_t conditional on knowing $(W_{t-1}, u_{t-1}, \zeta_t)$, (which is equivalent to knowing W_t), is normal with mean $V_{\xi\zeta}V_{\zeta\zeta}^{-1}\zeta_t$ and covariance matrix $V_{\xi\xi} - V_{\xi\zeta}V_{\zeta\zeta}^{-1}V_{\zeta\xi}$. These give (16.1)–(16.3). \square

16.2 Certainty equivalence

We say that a quantity a is *policy-independent* if $E_\pi(a | W_0)$ is independent of π .

Theorem 16.2. *Suppose LQG model assumptions hold. Then (i) the value function is of the form*

$$F(W_t) = \hat{x}_t^\top \Pi_t \hat{x}_t + \dots \quad (16.4)$$

where \hat{x}_t is the linear least squares estimate of x_t whose evolution is determined by the Kalman filter in Theorem 16.1 and ‘ $+\dots$ ’ indicates terms that are policy independent; (ii) the optimal control is given by

$$u_t = K_t \hat{x}_t,$$

where Π_t and K_t are the same matrices as in the full information case of Theorem 13.2.

It is important to grasp the remarkable fact that (ii) asserts: *the optimal control u_t is exactly the same as it would be if all unknowns were known and took values equal to their linear least square estimates (equivalently, their conditional means) based upon observations up to time t .* This is the idea known as **certainty equivalence**. As we have seen in the previous section, the distribution of the estimation error $\hat{x}_t - x_t$ does not depend on U_{t-1} . The fact that the problems of optimal estimation and optimal control can be decoupled in this way is known as the **separation principle**.

Proof. The proof is by backward induction. Suppose (16.4) holds at t . Recall that

$$\hat{x}_t = A\hat{x}_{t-1} + Bu_{t-1} + H_t\zeta_t, \quad \Delta_{t-1} = \hat{x}_{t-1} - x_{t-1}.$$

Then with a quadratic cost of the form $c(x, u) = x^\top Rx + 2u^\top Sx + u^\top Qu$, we have

$$\begin{aligned}
F(W_{t-1}) &= \min_{u_{t-1}} E [c(x_{t-1}, u_{t-1}) + \hat{x}_t^\top \Pi_t \hat{x}_t + \dots \mid W_{t-1}, u_{t-1}] \\
&= \min_{u_{t-1}} E \left[c(\hat{x}_{t-1} - \Delta_{t-1}, u_{t-1}) \right. \\
&\quad + (A\hat{x}_{t-1} + Bu_{t-1} + H_t \zeta_t)^\top \Pi_t (A\hat{x}_{t-1} + Bu_{t-1} + H_t \zeta_t) \\
&\quad \left. + \dots \mid W_{t-1}, u_{t-1} \right] \\
&= \min_{u_{t-1}} [c(\hat{x}_{t-1}, u_{t-1}) + (A\hat{x}_{t-1} + Bu_{t-1})^\top \Pi_t (A\hat{x}_{t-1} + Bu_{t-1})] + \dots,
\end{aligned} \tag{16.5}$$

where we use the fact that, conditional on W_{t-1}, u_{t-1} , the quantities Δ_{t-1} and ζ_t have zero means and are policy independent. So when we evaluate (16.5) the expectations of all terms which are linear in these quantities are zero, like $E[u_{t-1}^\top S \Delta_{t-1}]$, and the expectations of all terms which are quadratic in these quantities, like $E[\Delta_{t-1}^\top R \Delta_{t-1}]$, are policy independent (and so may be included as part of $+\dots$). \square

16.3 The Hamilton-Jacobi-Bellman equation

In continuous time the plant equation is,

$$\dot{x} = a(x, u, t).$$

Consider a discounted cost of

$$\mathbf{C} = \int_0^h e^{-\alpha t} c(x, u, t) dt + e^{-\alpha h} \mathbf{C}(x(h), h).$$

The discount factor over δ is $e^{-\alpha\delta} = 1 - \alpha\delta + o(\delta)$. So the optimality equation is,

$$F(x, t) = \inf_u [c(x, u, t)\delta + e^{-\alpha\delta} F(x + a(x, u, t)\delta, t + \delta) + o(\delta)].$$

By considering the term of order δ in the Taylor series expansion we obtain,

$$\inf_u \left[c(x, u, t) - \alpha F + \frac{\partial F}{\partial t} + \frac{\partial F}{\partial x} a(x, u, t) \right] = 0, \quad t < h, \tag{16.6}$$

with $F(x, h) = \mathbf{C}(x, h)$. In the undiscounted case, we simply put $\alpha = 0$. Notice that in (17.9) we have $\alpha = 0$ and the term of $\frac{\partial F}{\partial t}$ disappears because $h = \infty$.

Equation (16.6) is called the **Hamilton-Jacobi-Bellman equation** (HJB). Its heuristic derivation we have given above is justified by the following theorem. It can be viewed as the equivalent, in continuous time, of the backwards induction that we use in discrete time to verify that a policy is optimal because it satisfies the the dynamic programming equation.

Theorem 16.3. *Suppose a policy π , using a control u , has a value function F which satisfies the HJB equation (16.6) for all values of x and t . Then π is optimal.*

Proof. Consider any other policy, using control v , say. Then along the trajectory defined by $\dot{x} = a(x, v, t)$ we have

$$\begin{aligned} -\frac{d}{dt}e^{-\alpha t}F(x, t) &= e^{-\alpha t} \left[c(x, v, t) - \left(c(x, v, t) - \alpha F + \frac{\partial F}{\partial t} + \frac{\partial F}{\partial x} a(x, v, t) \right) \right] \\ &\leq e^{-\alpha t} c(x, v, t). \end{aligned}$$

The inequality is because the term round brackets is nonnegative. Integrating this inequality along the v path, from $x(0)$ to $x(h)$, gives

$$F(x(0), 0) - e^{-\alpha h} \mathbf{C}(x(h), h) \leq \int_{t=0}^h e^{-\alpha t} c(x, v, t) dt.$$

Thus the v path incurs a cost of at least $F(x(0), 0)$, and hence π is optimal. \square

16.4 Example: LQ regulation

The undiscounted continuous time DP equation for the LQ regulation problem is

$$0 = \inf_u [x^\top R x + u^\top Q u + F_t + F_x^\top (A x + B u)].$$

Suppose we try a solution of the form $F(x, t) = x^\top \Pi(t)x$, where $\Pi(t)$ is a symmetric matrix. Then $F_x = 2\Pi(t)x$ and the optimizing u is $u = -\frac{1}{2}Q^{-1}B^\top F_x = -Q^{-1}B^\top \Pi(t)x$. Therefore the DP equation is satisfied with this u if

$$0 = x^\top \left[R + \Pi A + A^\top \Pi - \Pi B Q^{-1} B^\top \Pi + \frac{d\Pi}{dt} \right] x,$$

where we use the fact that $2x^\top \Pi A x = x^\top \Pi A x + x^\top A^\top \Pi x$. This must hold for all x . So we have a solution to the HJB equation if $\Pi(t)$ satisfies the Riccati differential equation

$$R + \Pi A + A^\top \Pi - \Pi B Q^{-1} B^\top \Pi + \frac{d\Pi}{dt} = 0,$$

with a given boundary value for $\Pi(h)$.

16.5 Example: harvesting fish

A fish population of size x obeys the plant equation,

$$\dot{x} = a(x, u) = \begin{cases} a(x) - u & x > 0, \\ a(x) & x = 0. \end{cases}$$

The function $a(x)$ reflects the facts that the population can grow when it is small, but is subject to environmental limitations when it is large. It is desired to maximize the discounted total harvest $\int_0^T u e^{-\alpha t} dt$, subject to $0 \leq u \leq u_{\max}$, where u_{\max} is the greatest possible fishing rate.

Solution. The DP equation (with discounting) is

$$\sup_u \left[u - \alpha F + \frac{\partial F}{\partial t} + \frac{\partial F}{\partial x} [a(x) - u] \right] = 0, \quad t < T.$$

Since u occurs linearly with the maximization we again have a bang-bang optimal control, of the form

$$u = \begin{bmatrix} 0 \\ \text{undetermined} \\ u_{\max} \end{bmatrix} \text{ for } F_x \begin{bmatrix} > \\ = \\ < \end{bmatrix} 1.$$

Suppose $F(x, t) \rightarrow F(x)$ as $T \rightarrow \infty$, and $\partial F / \partial t \rightarrow 0$. Then

$$\sup_u \left[u - \alpha F + \frac{\partial F}{\partial x} [a(x) - u] \right] = 0. \quad (16.7)$$

Let us make a guess that $F(x)$ is concave, and then deduce that

$$u = \begin{bmatrix} 0 \\ \text{undetermined, but effectively } a(\bar{x}) \\ u_{\max} \end{bmatrix} \text{ for } x \begin{bmatrix} < \\ = \\ > \end{bmatrix} \bar{x}. \quad (16.8)$$

Clearly, \bar{x} is the operating point. We suppose

$$\dot{x} = \begin{cases} a(x) > 0, & x < \bar{x} \\ a(x) - u_{\max} < 0, & x > \bar{x}. \end{cases}$$

We say that there is **chattering** about the point \bar{x} , in the sense that u will switch between its maximum and minimum values either side of \bar{x} , effectively taking the value $a(\bar{x})$ at \bar{x} . To determine \bar{x} we note that

$$F(\bar{x}) = \int_0^\infty e^{-\alpha t} a(\bar{x}) dt = a(\bar{x}) / \alpha. \quad (16.9)$$

So from (16.7) and (16.9) we have

$$F_x(x) = \frac{\alpha F(x) - u(x)}{a(x) - u(x)} \rightarrow 1 \text{ as } x \nearrow \bar{x} \text{ or } x \searrow \bar{x}. \quad (16.10)$$

For F to be concave, F_{xx} must be negative if it exists. So we must have

$$\begin{aligned} F_{xx} &= \frac{\alpha F_x}{a(x) - u} - \left(\frac{\alpha F - u}{a(x) - u} \right) \left(\frac{a'(x)}{a(x) - u} \right) \\ &= \left(\frac{\alpha F - u}{a(x) - u} \right) \left(\frac{\alpha - a'(x)}{a(x) - u} \right) \\ &\simeq \frac{\alpha - a'(x)}{a(x) - u(x)} \end{aligned}$$

where the last line follows because (16.10) holds in a neighbourhood of \bar{x} . It is required that F_{xx} be negative. But the denominator changes sign at \bar{x} , so the numerator must do so also, and therefore we must have $a'(\bar{x}) = \alpha$.

We now have the complete solution. The control in (16.8) has a value function F which satisfies the HJB equation.

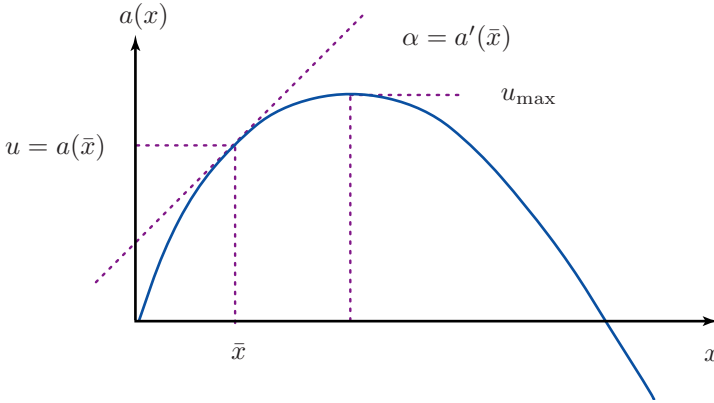


Figure 2: Growth rate $a(x)$ subject to environment pressures

Notice that we sacrifice long term yield for immediate return. If the initial population is greater than \bar{x} then the optimal policy is to fish at rate u_{\max} until we reach \bar{x} and then fish at rate $u = a(\bar{x})$. As $\alpha \nearrow a'(0)$, $\bar{x} \searrow 0$. If $\alpha \geq a'(0)$ then it is optimal to wipe out the entire fish stock.

Finally, it would be good to verify that $F(x)$ is concave, as we conjectured from the start. To see this, suppose $x > \bar{x}$. Then

$$\begin{aligned} F(x) &= \int_0^T u_{\max} e^{-\alpha t} dt + \int_T^\infty a(\bar{x}) e^{-\alpha t} dt \\ &= a(\bar{x})/\alpha + (u_{\max} - a(\bar{x}))(1 - e^{-\alpha T})/\alpha \end{aligned}$$

where $T = T(x)$ is the time taken for the fish population to decline from x to \bar{x} , when $\dot{x} = a(x) - u_{\max}$. Now

$$\begin{aligned} T(x) &= \delta + T(x + (a(x) - u_{\max})\delta) \implies 0 = 1 + (a(x) - u_{\max})T'(x) \\ &\implies T'(x) = 1/(u_{\max} - a(x)) \end{aligned}$$

So $F''(x)$ has the same sign as that of

$$\frac{d^2}{dx^2} (1 - e^{-\alpha T}) = -\frac{\alpha e^{-\alpha T} (\alpha - a'(x))}{(u_{\max} - a(x))^2},$$

which is negative, as required, since $\alpha = a'(\bar{x}) \geq a'(x)$, when $x > \bar{x}$. The case $x < \bar{x}$ is similar.

17 Pontryagin's Maximum Principle

Pontryagin's maximum principle. Optimization of consumption. Parking a rocket car. Adjoint variables as Lagrange multipliers.

17.1 Example: optimization of consumption

Suppose that given $x(0)$, κ and T , all positive, we wish to choose $u(t)$ to maximize

$$\int_0^T \log u(t) dt + \kappa \log x(T), \quad \text{subject to } \dot{x}(t) = ax(t) - u(t), \quad 0 \leq t \leq T.$$

Solution. Try using a Lagrange multiplier $\lambda(t)$ for the constraint $\dot{x}(t) = ax(t) - u(t)$. The Lagrangian is

$$L = \kappa \log x(T) + \int_0^T [\log u - \lambda(\dot{x} - (ax - u))] dt$$

Now use integration by parts, and define $H(x, u, \lambda) = \log u + \lambda(ax - u)$.

$$\begin{aligned} L &= \kappa \log x(T) - \lambda(t)x(t) \Big|_0^T + \int_0^T [\log u + \dot{\lambda}x + \lambda(ax - u)] dt \\ &= [\kappa \log x(T) - \lambda(T)x(T)] + \lambda(0)x(0) + \int_0^T [\dot{\lambda}x + H(x, u, \lambda)] dt \end{aligned}$$

For L to be stationary with respect to both $x(t)$ and $u(t)$, at every point within the integrand, we need

$$\begin{aligned} \dot{\lambda} + \frac{\partial}{\partial x} H(x, u, \lambda) &= 0 \implies \dot{\lambda} = -a\lambda \\ \frac{\partial}{\partial u} H(x, u, \lambda) &= 0 \implies u = 1/\lambda, \end{aligned}$$

and so $\lambda(t) = \lambda(0)e^{-at}$, $u(t) = \lambda(0)^{-1}e^{at}$ and $\dot{x}(t) = ax(t) - \lambda(0)^{-1}e^{at}$.

If the value of $x(T)$ is prescribed (and $< e^{aT}x(0)$ so u need not be negative), then we can solve this differential equation for x , choosing $\lambda(0)$ so that $x(t)$ takes prescribed values $x(0)$ and $x(T)$ at $t = 0$ and T . We get (after some algebra)

$$u(t) = \left(\frac{x(0) - x(T)e^{-aT}}{(T-t)x(0) - x(T)e^{-aT}} \right) x(t).$$

If the value of $x(T)$ is free, then stationarity of L with respect to $x(T)$ requires $\kappa/x(T) - \lambda(T) = 0$ which (after some algebra) implies $\lambda(0) = (\kappa + T)/x(0)$ and

$$u(t) = \frac{1}{T + \kappa} x(0) e^{at} = \frac{1}{(T-t) + \kappa} x(t). \quad (17.1)$$

If $a > 1/(\kappa + T)$ the trajectory is one in which $x(t)$ is initially increasing and then decreasing; otherwise $x(t)$ is decreasing. The optimal 'inheritance' left at T is

$$x(T) = \frac{\kappa}{\kappa + T} x(0) e^{aT}.$$

17.2 Heuristic derivation of Pontryagin's maximum principle

Pontryagin's maximum principle (PMP) states a *necessary condition that must hold on an optimal trajectory*. It is a calculation for a *fixed* initial value of the state, $x(0)$. In comparison, the dynamic programming approach is a calculation for a general initial value of the state. Thus, when PMP is useful, it finds an open-loop prescription of the optimal control, whereas dynamic programming is useful for finding a closed-loop prescription. PMP can be used as both a computational and analytic technique (and in the second case can solve the problem for general initial value.)

We begin by considering a time-homogeneous formulation, with plant equation $\dot{x} = a(x, u)$ and instantaneous cost $c(x, u)$. The trajectory is to be controlled until it reaches some stopping set S , where there is a terminal cost $K(x)$. As in (16.6) the value function $F(x)$ obeys the a dynamic programming equation (without discounting)

$$\inf_{u \in \mathcal{U}} \left[c(x, u) + \frac{\partial F}{\partial x} a(x, u) \right] = 0, \quad x \notin S, \quad (17.2)$$

with terminal condition

$$F(x) = K(x), \quad x \in S. \quad (17.3)$$

Define the **adjoint variable**

$$\lambda = -F_x. \quad (17.4)$$

This is a column n -vector, and is to be regarded as a function of time as the state moves along the optimal trajectory. The proof that F_x exists in the required sense is actually a tricky technical matter. We also define the **Hamiltonian**

$$H(x, u, \lambda) = \lambda^\top a(x, u) - c(x, u), \quad (17.5)$$

a scalar, defined at each point of the path as a function of the current x , u and λ .

Theorem 17.1. (PMP) *Suppose $u(t)$ and $x(t)$ represent the optimal control and state trajectory. Then there exists an adjoint trajectory $\lambda(t)$ such that together $u(t)$, $x(t)$ and $\lambda(t)$ satisfy*

$$\dot{x} = H_\lambda, \quad [= a(x, u)] \quad (17.6)$$

$$\dot{\lambda} = -H_x, \quad [= -\lambda^\top a_x + c_x] \quad (17.7)$$

and for all t , $0 \leq t \leq T$, and all feasible controls v ,

$$H(x(t), v, \lambda(t)) \leq H(x(t), u(t), \lambda(t)), \quad (17.8)$$

i.e. the optimal control $u(t)$ is the value of v maximizing $H((x(t), v, \lambda(t))$.

'Proof.' Our heuristic proof is based upon the DP equation; this is the most direct and enlightening way to derive conclusions that may be expected to hold in general.

Assertion (17.6) is immediate, and (17.8) follows from the fact that the minimizing value of u in (17.2) is optimal. Assuming u is the optimal control we have from (17.2) in incremental form as

$$F(x, t) = c(x, u)\delta + F(x + a(x, u)\delta, t + \delta) + o(\delta).$$

Now use the chain rule to differentiate with respect to x_i and this yields

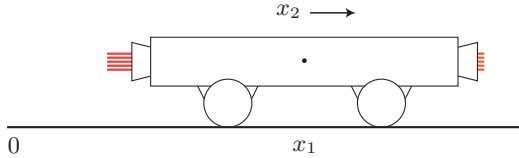
$$\begin{aligned} \frac{d}{dx_i} F(x, t) &= \delta \frac{d}{dx_i} c(x, u) + \sum_j \frac{\partial}{\partial x_j} F(x + a(x, u)\delta, t + \delta) \frac{d}{dx_i} (x_j + a_j(x, u)\delta) \\ \implies -\lambda_i(t) &= \delta \frac{dc}{dx_i} - \lambda_i(t + \delta) - \delta \sum_j \lambda_j(t + \delta) \frac{da_j}{dx_i} + o(\delta) \\ \implies \frac{d}{dt} \lambda_i(t) &= \frac{dc}{dx_i} - \sum_j \lambda_j(t) \frac{da_j}{dx_i} \end{aligned}$$

which is (17.7). □

Notice that (17.6) and (17.7) each give n equations. Condition (17.8) gives a further m equations (since it requires stationarity with respect to variation of the m components of u .) So in principle these equations, if nonsingular, are sufficient to determine the $2n + m$ functions $u(t)$, $x(t)$ and $\lambda(t)$.

17.3 Example: parking a rocket car

A rocket car has engines at both ends. Initial position and velocity are $x_1(0)$ and $x_2(0)$.



By firing the rockets (causing acceleration of u in the forward or reverse direction) we wish to park the car in minimum time, i.e. minimize T such that $x_1(T) = x_2(T) = 0$. The dynamics are $\dot{x}_1 = x_2$ and $\dot{x}_2 = u$, where u is constrained by $|u| \leq 1$.

Let $F(x)$ be minimum time that is required to park the rocket car. Then

$$F(x_1, x_2) = \min_{-1 \leq u \leq 1} \left\{ \delta + F(x_1 + x_2\delta, x_2 + u\delta) \right\}.$$

By making a Taylor expansion and then letting $\delta \rightarrow 0$ we find the HJB equation:

$$0 = \min_{-1 \leq u \leq 1} \left\{ 1 + \frac{\partial F}{\partial x_1} x_2 + \frac{\partial F}{\partial x_2} u \right\} \quad (17.9)$$

with boundary condition $F(0,0) = 0$. We can see that the optimal control will be a **bang-bang control** with $u = -\text{sign}(\frac{\partial F}{\partial x_2})$ and so F satisfies

$$0 = 1 + \frac{\partial F}{\partial x_1} x_2 - \left| \frac{\partial F}{\partial x_2} \right|.$$

Now let us tackle the same problem using PMP. We wish to minimize

$$\mathbf{C} = \int_0^T 1 dt$$

where T is the first time at which $x = (0,0)$. For dynamics if $\dot{x}_1 = x_2$, $\dot{x}_2 = u$, $|u| \leq 1$, the Hamiltonian is

$$H = \lambda_1 x_2 + \lambda_2 u - 1,$$

which is maximized by $u = \text{sign}(\lambda_2)$. The adjoint variables satisfy $\dot{\lambda}_i = -\partial H / \partial x_i$, so

$$\dot{\lambda}_1 = 0, \quad \dot{\lambda}_2 = -\lambda_1. \quad (17.10)$$

Suppose that at termination $\lambda_1 = \alpha$, $\lambda_2 = \beta$. Then in terms of time to go we can compute

$$\lambda_1 = \alpha, \quad \lambda_2 = \beta + \alpha s.$$

These reveal the form of the solution: there is at most one change of sign of λ_2 on the optimal path; u is maximal in one direction and then possibly maximal in the other.

From (17.2) or (17.9) we see that the maximized value of H must be 0. So at termination (when $x_2 = 0$), we conclude that we must have $|\beta| = 1$. We now consider the case $\beta = 1$. The case $\beta = -1$ is similar.

If $\beta = 1$, $\alpha \geq 0$ then $\lambda_2 = 1 + \alpha s \geq 0$ for all $s \geq 0$ and

$$u = 1, \quad x_2 = -s, \quad x_1 = s^2/2.$$

In this case the optimal trajectory lies on the parabola $x_1 = x_2^2/2$, $x_1 \geq 0$, $x_2 \leq 0$. This is half of the **switching locus** $x_1 = \pm x_2^2/2$ (shown dotted in Figure 3). Notice that the path is sensitive to the initial conditions. Consider a and b , just either side of the switching locus. From a we take first $u = 1$ then $u = -1$. From b we first take $u = -1$, then $u = 1$.

If $\beta = 1$, $\alpha < 0$ then $u = -1$ or $u = 1$ as the time to go is greater or less than $s_0 = 1/|\alpha|$. In this case,

$$\begin{aligned} u = -1, \quad x_2 = (s - 2s_0), \quad x_1 = 2s_0s - \frac{1}{2}s^2 - s_0^2, & \quad s \geq s_0, \\ u = 1, \quad x_2 = -s, \quad x_1 = \frac{1}{2}s^2, & \quad s \leq s_0. \end{aligned}$$

The control rule expressed as a function of s is open-loop, but in terms of (x_1, x_2) and the switching locus, it is closed-loop.

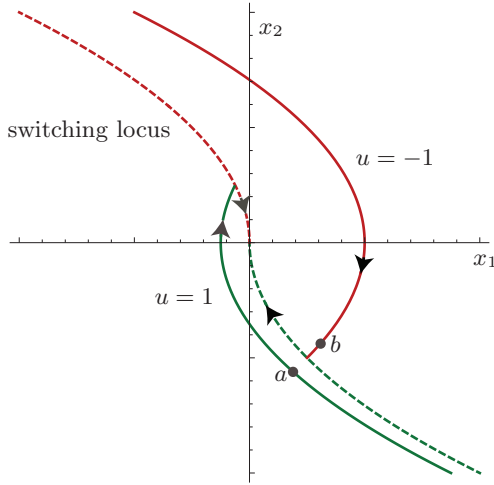


Figure 3: Optimal trajectories for parking a rocket car. Notice that the trajectories starting from two nearby points, a and b , are qualitatively different.

17.4 Adjoint variables as Lagrange multipliers

We have already seen in §17.1 that it is possible to think of $\lambda(t)$ as a Lagrange multiplier for the constraint $\dot{x} = a(x, u)$ (at time t). Consider the Lagrangian

$$L = -K(x(T)) + \int_0^T [-c - \lambda^\top (\dot{x} - a)] dt.$$

This is to be maximized over (x, u, λ) paths having the property that $x(t)$ first enters the set S at time T . We integrate $\lambda^\top \dot{x}$ by parts to obtain

$$L = -K(x(T)) - \lambda(T)^\top x(T) + \lambda(0)^\top x(0) + \int_0^T [\dot{\lambda}^\top x + \lambda^\top a - c] dt.$$

We now think about varying both $x(t)$ and $u(t)$, but without regard to the constraint $\dot{x} = a(x, u)$. The quantity within the integral must be stationary with respect to $x = x(t)$ and hence $\dot{\lambda} + \lambda^\top a_x - c_x = 0 \implies \dot{\lambda} = -H_x$, i.e. (17.7).

If $x(T)$ is unconstrained then the Lagrangian must also be stationary with respect to small variations in $x(T)$ that are in a direction σ such that $x(T) + \epsilon\sigma$ is in the stopping set (or within $o(\epsilon)$ of it), and this gives $(K_x(x(T)) + \lambda(T))^\top \sigma = 0$, i.e. the so-called **transversality conditions**, which we will say more about in (18.1).

It is good to have this alternative viewpoint, but it is informal and less easy to rigourise than the ‘proofs’ of in §17.2, and §18.1

18 Using Pontryagin's Maximum Principle

Transversality conditions. Examples with Pontryagin's maximum principle.

18.1 Transversality conditions

In (17.2) we see that H must be maximized to 0. We can make this a generally valid assertion, and also say some things about the terminal value of $\lambda(T)$ (the so-called **transversality conditions**.)

Theorem 18.1. (i) $H = 0$ on the optimal path. (ii) The terminal condition

$$(\lambda + K_x)^\top \sigma = 0 \tag{18.1}$$

holds at the terminal x for all σ such that $x + \epsilon\sigma$ is within $o(\epsilon)$ of the termination point of a possible optimal trajectory for all sufficiently small positive ϵ .

'Proof.' Assertion (i) follows from (17.2). To see (ii), suppose that x is a point at which the optimal trajectory first enters S . Then $x \in S$ and so $F(x) = K(x)$. Suppose $x + \epsilon\sigma + o(\epsilon) \in S$. Then

$$\begin{aligned} 0 &= F(x + \epsilon\sigma + o(\epsilon)) - K(x + \epsilon\sigma + o(\epsilon)) \\ &= F(x) - K(x) + (F_x(x) - K_x(x))^\top \sigma \epsilon + o(\epsilon) \end{aligned}$$

Together with $F(x) = K(x)$ this gives $(F_x - K_x)^\top \sigma = 0$. Since $\lambda = -F_x$ we get $(\lambda + K_x)^\top \sigma = 0$. \square

18.2 Example: use of transversality conditions

Suppose $\dot{x}_1 = x_2$, $\dot{x}_2 = u$, $x(0) = (0, 0)$, u is unconstrained, and we wish to minimize

$$\mathbf{C} = -x_1(1) + \int_0^1 \frac{1}{2}u(t)^2 dt.$$

Here $K(x) = -x_1(1)$. The Hamiltonian is

$$H(x, u, \lambda) = \lambda_1 x_2 + \lambda_2 u - \frac{1}{2}u^2,$$

which is maximized at $u(t) = \lambda_2(t)$. Now $\dot{\lambda}_i = -\partial H/\partial x_i$ gives

$$\dot{\lambda}_1 = 0, \quad \dot{\lambda}_2 = -\lambda_1.$$

The terminal x is unconstrained so in the transversality condition of $(\lambda + K_x)^\top \sigma = 0$, σ is arbitrary and so we also have

$$\lambda_1(1) - 1 = 0, \quad \lambda_2(1) = 0.$$

Thus the solution must be $\lambda_1(t) = 1$ and $\lambda_2(t) = 1 - t$. The optimal control is $u(t) = 1 - t$.

Note that there is often more than one way to set up a control problem. In this problem, we might have taken $K = 0$, but included a cost of $-\int_0^1 x_2 dt = -x_1(1) + x_1(0)$.

18.3 Example: insects as optimizers

A colony of insects consists of workers and queens, of numbers $w(t)$ and $q(t)$ at time t . If a time-dependent proportion $u(t)$ of the colony's effort is put into producing workers, ($0 \leq u(t) \leq 1$), then w, q obey the equations

$$\dot{w} = auw - bw, \quad \dot{q} = c(1 - u)w,$$

where a, b, c are constants, with $a > b$. The function u is to be chosen to maximize the number of queens at the end of the season. Show that the optimal policy is to produce only workers up to some moment, and produce only queens thereafter.

Solution. In this problem the Hamiltonian is

$$H = \lambda_1(auw - bw) + \lambda_2c(1 - u)w$$

and $K(w, q) = -q$. The adjoint equations and transversality conditions give

$$\begin{aligned} -\dot{\lambda}_1 &= H_w = \lambda_1(au - b) + \lambda_2c(1 - u) & \lambda_1(T) &= -K_w = 0 \\ -\dot{\lambda}_2 &= H_q = 0 & \lambda_2(T) &= -K_q = 1 \end{aligned} ,$$

So $\lambda_2(t) = 1$ for all t . Since H is maximized by u so

$$u = \begin{cases} 0 & \text{if } \Delta(t) := \lambda_1 a - c < 0 \\ 1 & > 0 \end{cases} .$$

Since $\Delta(T) = -c$, we must have $u(T) = 0$. If t is a little less than T , λ_1 is small and $u = 0$ so the equation for λ_1 is

$$\dot{\lambda}_1 = \lambda_1 b - c. \tag{18.2}$$

As long as λ_1 is small, $\dot{\lambda}_1 < 0$. Therefore as the *remaining time* s increases, $\lambda_1(s)$ increases, until such point that $\Delta(t) = \lambda_1 a - c \geq 0$. The optimal control becomes $u = 1$ and then $\dot{\lambda}_1 = -\lambda_1(a - b) < 0$, which implies that $\lambda_1(s)$ continues to increase as s increases, right back to the start. So there is no further switch in u .

The point at which the single switch occurs is found by integrating (18.2) from t to T , to give $\lambda_1(t) = (c/b)(1 - e^{-(T-t)b})$ and so the switch occurs where $\lambda_1 a - c = 0$, i.e. $(a/b)(1 - e^{-(T-t)b}) = 1$, or

$$t_{\text{switch}} = T + (1/b) \log(1 - b/a).$$

Experimental evidence suggests that social insects do closely follow this policy and adopt a switch time that is nearly optimal for their natural environment.

18.4 Problems in which time appears explicitly

Thus far, $c(\cdot)$, $a(\cdot)$ and $K(\cdot)$ have been function of (x, u) , but not t . Sometimes we wish to solve problems in t appears, such as when $\dot{x} = a(x, u, t)$. We can cope with this

generalization by the simple mechanism of introducing a new variable that equates to time. Let $x_0 = t$, with $\dot{x}_0 = a_0 = 1$.

Having been augmented by this variable, the Hamiltonian gains a term and becomes

$$\tilde{H} = \lambda_0 a_0 + H = \lambda_0 a_0 + \sum_{i=1}^n \lambda_i a_i - c$$

where $\lambda_0 = -F_t$ and $a_0 = 1$. Theorem 18.1 says that \tilde{H} must be maximized to 0. Equivalently, on the optimal trajectory,

$$H(x, u, \lambda) = \sum_{i=1}^n \lambda_i a_i - c \text{ must be maximized to } -\lambda_0.$$

Theorem 17.1 still holds. However, to (17.7) we can now add

$$\dot{\lambda}_0 = -H_t = c_t - \lambda a_t, \tag{18.3}$$

and transversality condition

$$(\lambda + K_x)^\top \sigma + (\lambda_0 + K_t) \tau = 0, \tag{18.4}$$

which must hold at the termination point (x, t) if $(x + \epsilon \sigma, t + \epsilon \tau)$ is within $o(\epsilon)$ of the termination point of an optimal trajectory for all small enough positive ϵ .

We can now understand what to do with various types of time-dependency and terminal conditions on $x(T)$ and/or T . For example, we can draw the following inferences.

- (i) If K is time-independent (so $K_t = 0$) and the terminal time T is unconstrained (so τ is arbitrary) then the transversality condition implies $\lambda_0(T) = 0$. Since H is always maximized to $-\lambda_0(t)$ it must be maximized to 0 at T .
- (ii) If a, c are only functions of (x, u) then $\dot{\lambda}_0 = c_t - \lambda^\top a_t = 0$, and so $\lambda_0(t)$ is constant on the optimal trajectory. Since H is always maximized to $-\lambda_0(t)$ it must be maximized to a constant on the optimal trajectory.
- (iii) If both (i) and (ii) are true then H is maximized to 0 along the entire optimal trajectory. We had this in the problem of parking in minimal time, §17.3.

18.5 Example: monopolist

Miss Prout holds the entire remaining stock of Cambridge elderberry wine for the vintage year 1959. If she releases it at rate u (in continuous time) she realises a unit price $p(u) = (1 - u/2)$, for $0 \leq u \leq 2$ and $p(u) = 0$ for $u \geq 2$. She holds an amount x at time 0 and wishes to release it in a way that maximizes her total discounted return, $\int_0^T e^{-\alpha t} u p(u) dt$, (where T is unconstrained.)

Solution. Notice that t appears in the cost function. The plant equation is $\dot{x} = -u$ and the Hamiltonian is

$$H(x, u, \lambda) = e^{-\alpha t} u p(u) - \lambda u = e^{-\alpha t} u(1 - u/2) - \lambda u.$$

Note that $K = 0$. Maximizing with respect to u and using $\dot{\lambda} = -H_x$ gives

$$u = 1 - \lambda e^{\alpha t}, \quad \dot{\lambda} = 0, \quad t \geq 0,$$

so λ is constant. The terminal time is unconstrained so the transversality condition gives $\lambda_0(T) = -K_t|_{t=T} = 0$. Therefore, since we require H to be maximized to $-\lambda_0(T) = 0$ at T , we have $u(T) = 0$, and hence

$$\lambda = e^{-\alpha T}, \quad u = 1 - e^{-\alpha(T-t)}, \quad t \leq T,$$

where T is then the time at which all wine has been sold, and so

$$x = \int_0^T u dt = T - (1 - e^{-\alpha T}) / \alpha.$$

Thus u is implicitly a function of x , through T .

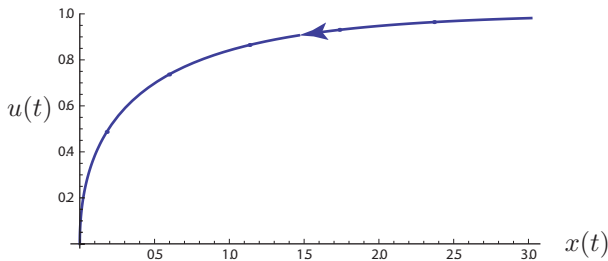


Figure 4: Trajectories of $x(t), u(t)$, for $\alpha = 1$.

The optimal value function is

$$F(x) = \int_0^T (u - u^2/2) e^{-\alpha t} dt = \frac{1}{2} \int_0^T (e^{-\alpha t} - e^{\alpha t - 2\alpha T}) dt = \frac{(1 - e^{-\alpha T})^2}{2\alpha}.$$

18.6 Example: neoclassical economic growth

Suppose x is the existing capital per worker and u is consumption of capital per worker. The plant equation is

$$\dot{x} = f(x) - \gamma x - u, \tag{18.5}$$

where $f(x)$ is production per worker (which depends on capital available to the worker), and $-\gamma x$ represents depreciation of capital. We wish to choose u to maximize

$$\int_{t=0}^T e^{-\alpha t} g(u) dt,$$

where $g(u)$ measures utility and T is prescribed.

Solution. This is really the same as the fish harvesting example in §16.5, with $a(x) = f(x) - \gamma x$. So let us take

$$\dot{x} = a(x) - u. \quad (18.6)$$

It is convenient to take

$$H = e^{-\alpha t} [g(u) + \lambda(a(x) - u)]$$

so including a discount factor in the definition of u , corresponding to expression of F in terms of present values. Here λ is a scalar. Then $g'(u) = \lambda$ (assuming the maximum is at a stationary point), and

$$\frac{d}{dt} (e^{-\alpha t} \lambda) = -H_x = -e^{-\alpha t} \lambda a'(x) \quad (18.7)$$

or

$$\dot{\lambda}(t) = (\alpha - a'(x))\lambda(t). \quad (18.8)$$

From $g'(u) = \lambda$ we have $g''(u)\dot{u} = \dot{\lambda}$ and hence from (18.8) we obtain

$$\dot{u} = \frac{1}{\sigma(u)} [a'(x) - \alpha], \quad (18.9)$$

where

$$\sigma(u) = -\frac{g''(u)}{g'(u)}$$

is the elasticity of marginal utility. Assuming g is strictly increasing and concave we have $\sigma > 0$. So (x, u) are determined by (18.6) and (18.9). An equilibrium solution at \bar{x}, \bar{u} is determined by

$$\bar{u} = a(\bar{x}) \quad a'(\bar{x}) = \alpha,$$

These give the balanced growth path; interestingly, it is independent of g .

This provides an example of so-called **turnpike theory**. For sufficiently large T the optimal trajectory will move from the initial $x(0)$ to within an arbitrary neighbourhood of the balanced growth path (the turnpike) and stay there for all but an arbitrarily small fraction of the time. As the terminal time becomes imminent the trajectory leaves the neighbourhood of the turnpike and heads for the terminal point $x(T) = 0$.

19 Controlled Diffusion Processes

Control problems in a continuous-time, continuous state space, stochastic setting.

19.1 The dynamic programming equation

The DP equation in incremental form is

$$F(x, t) = \inf_u \{c(x, u)\delta t + E[F(x(t + \delta t), t + \delta t) \mid x(t) = x, u(t) = u]\}.$$

If appropriate limits exist then this can be written in the limit $\delta t \downarrow 0$ as

$$\inf_u [c(x, u) + F_t(x, t) + \Lambda(u)F(x, t)] = 0.$$

Here $\Lambda(u)$ is the operator defined by

$$\Lambda(u)\phi(x) = \lim_{\delta t \downarrow 0} \left[\frac{E[\phi(x(t + \delta t)) \mid x(t) = x, u(t) = u] - \phi(x)}{\delta t} \right] \quad (19.1)$$

or

$$\Lambda(u)\phi(x) = \lim_{\delta t \downarrow 0} E \left[\frac{\phi(x(t + \delta t)) - \phi(x)}{\delta t} \mid x(t) = x, u(t) = u \right]$$

the conditional expectation of the ‘rate of change’ of $\phi(x)$ along the path. The operator Λ converts a scalar function of state, $\phi(x)$, to another such function, $\Lambda\phi(x)$. However, its action depends upon the control u , so we write it as $\Lambda(u)$. It is called the **infinitesimal generator** of the controlled Markov process. Equation (19.1) is equivalent to

$$E[\phi(x(t + \delta t)) \mid x(t) = x, u(t) = u] = \phi(x) + \Lambda(u)\phi(x)\delta t + o(\delta t).$$

This equation takes radically different forms depending upon whether the state space is discrete or continuous. Both are important.

If the state space is discrete we have the Markov jump process of Lecture 9. In this case $\Lambda(u)\phi(i) = \sum_j q_{ij}(u)[\phi(j) - \phi(i)]$. Now we turn to the case of continuous state space.

19.2 Diffusion processes and controlled diffusion processes

The **Wiener process** $\{B(t)\}$, is a scalar process for which $B(0) = 0$, the increments in B over disjoint time intervals are statistically independent and $B(t)$ is normally distributed with zero mean and variance t . (‘ B ’ stands for **Brownian motion**. It can be understood as a $\delta \rightarrow 0$ limit of a symmetric random walk in which steps $\pm\sqrt{\delta}$ are made at times $\delta, 2\delta, \dots$) The specification is internally consistent because, for example,

$$B(t) = B(t_1) + [B(t) - B(t_1)]$$

and for $0 \leq t_1 \leq t$ the two terms on the right-hand side are independent normal variables of zero mean and with variance t_1 and $t - t_1$ respectively.

If δB is the increment of B in a time interval of length δt then

$$E(\delta B) = 0, \quad E[(\delta B)^2] = \delta t, \quad E[(\delta B)^j] = o(\delta t), \quad \text{for } j > 2,$$

where the expectation is one conditional on the past of the process. Note that since

$$E[(\delta B/\delta t)^2] = O[(\delta t)^{-1}] \rightarrow \infty,$$

the formal derivative $\epsilon = dB/dt$ (continuous-time ‘white noise’) does not exist in a mean-square sense, but expectations such as

$$E \left[\left\{ \int \alpha(t)\epsilon(t)dt \right\}^2 \right] = E \left[\left\{ \int \alpha(t)dB(t) \right\}^2 \right] = \int \alpha(t)^2 dt$$

make sense if the integral is convergent.

Now consider a **stochastic differential equation**

$$\delta x = a(x, u)\delta t + g(x, u)\delta B,$$

which we shall write formally as

$$\dot{x} = a(x, u) + g(x, u)\epsilon.$$

This, as a Markov process, has an infinitesimal generator with action

$$\begin{aligned} \Lambda(u)\phi(x) &= \lim_{\delta t \downarrow 0} E \left[\frac{\phi(x(t + \delta t)) - \phi(x)}{\delta t} \middle| x(t) = x, u(t) = u \right] \\ &= \phi_x a + \frac{1}{2} \phi_{xx} g^2 \\ &= \phi_x a + \frac{1}{2} N \phi_{xx}, \end{aligned}$$

where $N(x, u) = g(x, u)^2$. So this **controlled diffusion process** has DP equation

$$\inf_u [c + F_t + F_x a + \frac{1}{2} N F_{xx}] = 0, \tag{19.2}$$

and in the vector case

$$\inf_u [c + F_t + F_x^\top a + \frac{1}{2} \text{tr}(N F_{xx})] = 0.$$

19.3 Example: noisy LQ regulation in continuous time

The dynamics are

$$\begin{aligned} \delta x &= Ax \delta t + Bu \delta t + N^{1/2} \delta B \\ \dot{x} &= Ax + Bu + N^{1/2} \epsilon. \end{aligned}$$

The dynamic programming equation is

$$\inf_u \left[x^\top R x + u^\top Q u + F_t + F_x^\top (A x + B u) + \frac{1}{2} \text{tr}(N F_{xx}) \right] = 0.$$

In analogy with the discrete and deterministic continuous cases that we have considered previously, we try a solution of the form,

$$F(x, t) = x^\top \Pi(t) x + \gamma(t).$$

This leads to the same Riccati equation as in Section 16.4,

$$0 = x^\top \left[R + \Pi A + A^\top \Pi - \Pi B Q^{-1} B^\top \Pi + \frac{d\Pi}{dt} \right] x,$$

and also, as in Section 13.3,

$$\frac{d\gamma}{dt} + \text{tr}(N \Pi(t)) = 0, \quad \text{giving} \quad \gamma(t) = \int_t^T \text{tr}(N \Pi(\tau)) d\tau.$$

19.4 Example: passage to a stopping set

Consider a problem of movement on the unit interval $0 \leq x \leq 1$ in continuous time, $\dot{x} = u + \epsilon$, where ϵ is white noise of power v . The process terminates at time T when x reaches one end or the other of the the interval. The cost is made up of an integral term $\frac{1}{2} \int_0^T (L + Q u^2) dt$, penalising both control and time spent, and a terminal cost which takes the value C_0 or C_1 according as termination takes place at 0 or 1.

Show that in the deterministic case $v = 0$ one should head straight for one of the termination points at a constant rate and that the value function $F(x)$ has a piecewise linear form, with possibly a discontinuity at one of the boundary points if that boundary point is the optimal target from no interior point of the interval.

Show, in the stochastic case, that the dynamic programming equation with the control value optimized out can be linearised by a transformation $F(x) = \alpha \log \phi(x)$ for a suitable constant α , and hence solve the problem.

Solution. In the deterministic case the optimality equation is

$$\inf_u \left[\frac{L + Q u^2}{2} + u \frac{\partial F}{\partial x} \right] = 0, \quad 0 < x < 1, \quad (19.3)$$

with boundary conditions $F(0) = C_0$, $F(1) = C_1$. If one goes (from x) for $x = 0$ at speed w one incurs a cost of $C_0 + (x/2w)(L + Q w^2)$ with a minimum over w value of $C_0 + x\sqrt{LQ}$. Indeed (19.3) is solved by

$$F(x) = \min \left[C_0 + x\sqrt{LQ}, C_1 + (1-x)\sqrt{LQ} \right].$$

The minimizing option determines the target and the optimal w is $\sqrt{L/Q}$.

In the stochastic case

$$\inf_u \left[\frac{L + Qu^2}{2} + u \frac{\partial F}{\partial x} + \frac{v}{2} \frac{\partial^2 F}{\partial x^2} \right] = 0.$$

So $u = -Q^{-1}F_x$ and

$$L - Q^{-1} \left(\frac{\partial F}{\partial x} \right)^2 + v \frac{\partial^2 F}{\partial x^2} = 0.$$

Make the transform $F(x) = -Qv \log \phi(x)$ so $\phi(x) = e^{-F(x)/Qv}$. Then the above equation simplifies to

$$Qv^2 \frac{\partial^2 \phi}{\partial x^2} - L\phi = 0,$$

with solution

$$\phi(x) = k_1 \exp\left(\frac{x}{v} \sqrt{L/Q}\right) + k_2 \exp\left(-\frac{x}{v} \sqrt{L/Q}\right).$$

We choose the constants k_1, k_2 to meet the two boundary conditions on F .

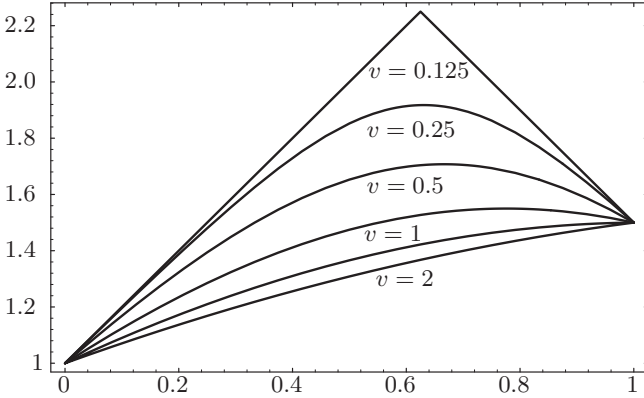


Figure 5: $F(x)$ against x for the passage to a stopping set

The figure shows the solution for $L = 1$, $Q = 4$, $C_0 = 1$, $C_1 = 1.5$ and $v = 0.125, 0.25, 0.5, 1, 2$ and the deterministic solution.

Notice that for these parameter choices the presence of noise actually reduces cost. This is because we are nearly indifferent as to which endpoint we hit, and L is small relative to Q . So it will be good to keep u small and let the noise do most of the work in bringing the state to an endpoint.

20 Risk-sensitive Optimal Control

A brief presentation of some ideas of Peter Whittle. Recapitulation of LQG model, certainty-equivalence and the maximum principle, but now with risk-sensitivity.

20.1 Whittle risk sensitivity

Consider a control problem in which the expected cost under policy π is $E_\pi C$. Whittle has proposed that one can model sensitivity to variability in the cost with the criterion

$$\gamma_\pi(\theta) = -\theta^{-1} \log [E_\pi (e^{-\theta C})].$$

Since the exponential function is convex, Jensen's inequality tells us that $E_\pi (e^{-\theta C}) \geq e^{-\theta E_\pi C}$, and so we see that $\gamma_\pi(\theta)$ is less, equal or greater than $E_\pi(C)$ as θ is positive, zero or negative. Observe also that

$$\gamma_\pi(\theta) = -\theta^{-1} \log \left[E_\pi \left(e^{-\theta E_\pi C - \theta(C - E_\pi C)} \right) \right] = E_\pi C - \frac{1}{2} \theta \text{var}_\pi(C) + \dots$$

So the variance of C enters as a first order term in θ . When θ is positive, zero or negative we are correspondingly risk-seeking, risk-neutral or risk-averse.

The LQG model with cost function $\gamma_\pi(C)$, and C a quadratic, is quite naturally given the acronym of LEQG (EQ meaning exponential of a quadratic).

20.2 The felicity of LEQG assumptions

Recall the discrete-time LQG state-structured model in which

$$\begin{aligned} x_t &= Ax_{t-1} + Bu_{t-1} + \epsilon_t \\ y_t &= Cx_{t-1} + \eta_t \end{aligned}$$

The noise terms are Gaussian, and initial information is $x_0 \sim N(\hat{x}_0, V_0)$. The cost function is

$$C(X, U) = \frac{1}{2} \sum_0^{t-1} (x_t^\top R x_t + u_t^\top Q u_t) + \frac{1}{2} x_h^\top \Pi_h x_h,$$

in which it is now convenient to supply an extra factor of $1/2$.

Under the LQG assumptions we have available to us many facts about least square estimation, maximum likelihood, partitioning of sums of squares, Gauss-Markov theorem, Shur-inverses, and so on. We now briefly summarise some of these and recall results from previous lectures.

Firstly, recall that it is frequently useful to partition a quadratic form Q , in vectors $x \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$, as

$$\begin{aligned} Q(x, u) &= \frac{1}{2} \begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} \Pi_{xx} & \Pi_{xu} \\ \Pi_{ux} & \Pi_{uu} \end{pmatrix} \begin{pmatrix} x \\ u \end{pmatrix} \\ &= \frac{1}{2} \left(u - \Pi_{uu}^{-1} \Pi_{ux} x \right)^\top \Pi_{uu} \left(u - \Pi_{uu}^{-1} \Pi_{ux} x \right) + \frac{1}{2} x^\top \left(\Pi_{xx} - \Pi_{xu} \Pi_{uu}^{-1} \Pi_{ux} \right) x. \quad (20.1) \end{aligned}$$

We are of course assuming that Π is symmetric and $\Pi_{uu} \succ 0$. Equation (20.1) shows that Q is minimized with respect to u by $\bar{u} = \Pi_{uu}^{-1}\Pi_{ux}x$, with the minimized value $Q(x, \bar{u}) = x^\top (\Pi_{xx} - \Pi_{xu}\Pi_{uu}^{-1}\Pi_{ux})x$. We have already used this in Lecture 13 to see that the optimal control for LQ regulation takes the form $u_t = K_t x_t$ and to derive the Riccati equation $\Pi_{t-1} = f\Pi_t$.

A second crucial fact is obtained by writing $h = u - \bar{u}$, and using (20.1) to give

$$\begin{aligned} \int_u e^{-Q(x,u)} du &= e^{-Q(x,\bar{u})} \int_u e^{-\frac{1}{2}h^\top \Pi_{uu} h} du = e^{-Q(x,\bar{u})} (\det(\Pi_{uu})(2\pi))^{m/2} \\ &= e^{-\inf_u Q(x,u)} (\det(\Pi_{uu})(2\pi))^{m/2}. \end{aligned} \quad (20.2)$$

The second integral can be evaluated by inspection if the reader remembers the form of the density function of a multivariate Gaussian $u \sim N(\bar{u}, \Pi_{uu}^{-1})$. The usefulness of (20.2) comes from noting that the dependence on x of the left-hand and right-hand sides is the same. This is also true in the more natural circumstance that the roles of u and x are interchanged, and we seek a control u which extremizes the expectation of $\exp(-\theta C)$ over a random x , for which $\exp(-Q(x, u))$ provides the x -dependent part of the Gaussian joint density function.

This is our third key fact: that exponentials of quadratics appear within the LEQG model in two ways. They appear in the cost $\exp(-\theta C)$, and also through the Gaussian joint density functions of unknown variables. If (x, u) are jointly Gaussian then their density function is proportional to $\exp(-Q(x, u))$, where Q is written in terms of the inverse of the covariance matrix V ,

$$Q(x, u) = \frac{1}{2} \begin{pmatrix} x \\ u \end{pmatrix}^\top \begin{pmatrix} V_{xx} & V_{xu} \\ V_{ux} & V_{uu} \end{pmatrix}^{-1} \begin{pmatrix} x \\ u \end{pmatrix}.$$

The conditional mean and covariance matrix of x given u are

$$\begin{aligned} E(x | u) &= V_{xu}V_{uu}^{-1}u \\ \text{cov}(x | u) &= V_{xx} - V_{xu}V_{uu}^{-1}V_{ux} = \tilde{S}. \end{aligned}$$

We had this as Lemma 15.4 when we derived the Kalman filter, (a recursive method of finding the parameters of the conditional distribution of x_t given W_t). We needed the fact that the conditional mean of ξ_t ($= x_t - A\hat{x}_t - Bu_{t-1}$) given η_t ($= y_t - C\hat{x}_t$) is a linear function of η_t and its conditional covariance matrix does not depend on the actual value of η_t .

As an aside, we remark that it is sometimes helpful to rewrite Q using the matrix $\tilde{S} = V_{xx} - V_{xu}V_{uu}^{-1}V_{ux}$, which is called the **Shur-complement** of V_{uu} in V . We might let $\Pi = V^{-1}$, in which case the following identity can be verified:

$$\begin{pmatrix} \Pi_{xx} & \Pi_{xu} \\ \Pi_{ux} & \Pi_{uu} \end{pmatrix} = \begin{pmatrix} \tilde{S}^{-1} & -\tilde{S}^{-1}V_{xu}V_{uu}^{-1} \\ V_{uu}^{-1}V_{ux}\tilde{S}^{-1} & V_{uu}^{-1} - V_{uu}^{-1}V_{ux}\tilde{S}^{-1}V_{xu}V_{uu}^{-1} \end{pmatrix}.$$

This helps to explain why the equations that specify optimal controls (Lecture 13) are so similar to those which specify optimal estimates (Lecture 15). We are equally content to compute quantities of interest from either blocks in V or in V^{-1} .

20.3 A risk-sensitive certainty-equivalence principle

Using the facts in §20.2 we can now address the problem of minimizing the risk-sensitive objective $\gamma_\pi(\theta)$. Let us present and then explain the following.

$$E_\pi [e^{-\theta C} | W_t] = \int e^{-\theta(C+\theta^{-1}D)} d\Box. \quad (20.3)$$

The cost C is quadratic in X, U . Given W_t , some constituent variables are known and others are unknown. Past controls are known, and future controls are not yet known. In the case of imperfect observation, X may never be known. The expectation is being taken under policy π with respect to all the unknown variables, denoted here by \Box . This expectation is computed by integrating $\exp(-\theta C)$ against the appropriate Gaussian density. Thus D (which Whittle calls the **discrepancy**) is essentially the quadratic form required to express the joint density function of all variables (both known and unknown) in terms of their covariance matrix.

Now we use the key fact (20.2) to write

$$\gamma_\pi(\theta) = -\theta^{-1} \log \left[e^{-\theta \inf_{\Box} (C+\theta^{-1}D)} \right] + \dots \quad (20.4)$$

where $+\dots$ denotes policy independent terms (which depend on $V_0, N, L, M, \Pi_h, R, Q$).

The nature of a **risk-sensitive certainly-equivalence principle** (RSCEP) is now clear. If at time t we are wishing to optimally choose u_t , then this can be viewed as a two-stage process. The first stage is to minimize the **stress**, defined as $S = C + \theta^{-1}D$, over all the unknown variables (the ensemble we are denoting as \Box), but excluding U . The LQG assumptions guarantee that these values will be linear functions of U . The second stage is to substitute these values into S and then minimize the resulting (quadratic) function of U . The Kalman filter expresses calculations of the first stage in a recursive manner. The second stage is simply a noiseless LQ regulation problem.

It is worth reflecting carefully on what (20.4) is saying. In the terms of our familiar notation, we are minimizing S over unknowns amongst X, Y and U . If we minimize over U first, then the optimizing U values will be linear function of X, Y , and we can subsequently minimize over these. By conducting the minimization in this sequence u_t can be chosen in full knowledge of W_t . Indeed, it looks as if we are allowed precognition when choosing u_t : potentially u_t can be a function of a future observable, such as y_τ for $\tau > t$ or even of an unobservable, such as x_0 . But now think about minimizing S with respect to X, Y before U . Now it is clear that U is no more than the sequence of controls that is optimal for the certainty-equivalent scenario that is obtained by replacing unknowns by their maximum likelihood estimates. This provides the risk-sensitive certainly-equivalence principle described above.

The special case as $\theta \rightarrow 0$ deserves comment. Now the term of $\theta^{-1}D$ overwhelmingly dominates S and so the first stage is to minimize D , which means replacing all unknowns with their maximum likelihood estimates, equivalently their means conditional on the values of the known variables. Then we solve the regulation problem under the assumption that these are the true values of the unknowns. In particular we

will estimate all future noise variables as 0. Having estimated all other unknowns, we must conclude that $u_t = K_t \hat{x}_t$, where K_t is the same matrix we would have used in the noiseless and full-information case. This is an alternative proof of Theorem 16.2, and in some respects a much cleaner and revealing one.

It is perhaps interesting to note that if one were to take all the following matrices which feature in the discrete-time LEQG model, namely $R, S, Q, L, M, N, V_0, \Pi_h$, and scale them by the same factor, then all estimates of unknowns do not change, and Π_t and V_t are simply scaled by the same factor.

We finish this section with a quote from Whittle, who writes

The RSCEP is in danger of being one of those pieces of work which is rejected for a time as odd and then is no sooner accepted than it is dismissed as trivial.

My personal view is that the reader who thinks at all deeply about (20.2) and (20.4) is unlikely to dismiss the RSCEP.

20.4 Large deviations

In the remainder of this chapter we will do no more than provide a few broad-brush hints about the way in which the theory of risk-sensitive control can be developed.

The RSCEP depends crucially on the fact that (20.2) turns the calculation of an expectation into a calculation of a path maximum likelihood. This is true only under LEQG assumptions. However, the theory of **large deviations** provides something very similar that can be applied more generally.

Suppose we have N copies of a stochastic path. Think, for example, of N i.i.d. random walks. The i th copy has $x_t^i = x_{t-1}^i + \epsilon_t^i$, $t = 0, \dots, h$, with $x_0^i = 0$. Let $Z_t = (1/N) \sum_{i=1}^N x_t^i$. The strong law of large numbers tells us that $\{Z_t, t = 0, \dots, h\}$ tends to the zero path almost surely as $N \rightarrow \infty$. We have already met this type of fluid model for a stochastic systems in §11.4.

Suppose we wish to evaluate $E[e^{-\theta C(Z)} \mid G]$, where cost C is a function of the path, and G is an event, like ‘*the path of Z ends near $Z_h = 0$, but $Z_t \geq 2$ for some t along the way*’. It is a part of large deviation theory to say that if such an unlikely event G occurs (which means a deviation from the zero path) then it is overwhelmingly most likely to occur in whatever way is most likely amongst possible ways that it could occur. In our example, this would mean making a straight-line path from 0 to 2 over $[0, h/2]$ and then another straight-line path from 2 to 0 over $[h/2, h]$. If this path is $\bar{\chi}$, then $E[e^{-\theta C(Z)} \mid G]$ is essentially $e^{-\theta C(\bar{\chi})}$. We have again turned the calculation of an expectation into a calculation of a path maximum likelihood. Moreover, if we introduce controls into our problem, then the calculation of the an optimal policy can be found by minimizing a stress of the form $C(\chi) + \theta^{-1} D(\chi)$. Once again we can compute this in stages, first with respect to unknown state and noise variables, and then with respect to controls. This means that a RSCEP can be valid outside the LEQG setting.

20.5 A risk-sensitive maximum principle

In §17.1 we used a Lagrange multiplier viewpoint to understand Pontryagin's maximum principle. We augmented the objective function by $\lambda^\top(\dot{x} - a(x, u))$ and then freely optimized over u and x .

Consider again the LEQG model, with full state observation, so $\dot{x} = Ax + Bu + \epsilon$. It is now plausible that we should add to the stress, $S = C + \theta^{-1}D$, a further term of $\lambda^\top(\dot{x} - Ax - Bu - \epsilon)$. The discrepancy part of this, $\theta^{-1}D$, already contains $(1/2)\theta^{-1}\epsilon^\top N\epsilon$, and so minimizing with respect to ϵ gives $\epsilon = \theta N\lambda$. After extremizing out ϵ we are left to extremize the path integral

$$\int_0^h [c(x, u) + \lambda^\top(\dot{x} - Ax - Bu) - \frac{1}{2}\theta\lambda^\top N\lambda] dt + C_h(x_h).$$

We can, as before integrate by parts to find a differential equation that λ must satisfy. We can also use the fact the dual variable λ must extremize the integrand in the opposite direction that in which we extremize with respect to x, u .

20.6 Example: risk-sensitive optimization of consumption

There is not space to continue with more theory. We simply conclude by returning to the example we studied in §17.1, namely the problem of maximizing consumption over a lifetime, but now with noise in the growth rate of capital:

$$\dot{x} = ax - u + \epsilon.$$

Let us take as a heuristic, which is motivated by the large deviations argument, that we should look at minimizing

$$\int_0^T [-\log u + \lambda(\dot{x} - (ax - u)) - \frac{1}{2}\theta N\lambda^2] dt - \log x(T).$$

If $\theta > 0$ we minimize with respect to x, u , and maximize with respect to λ . Thus $u = -1/\lambda$, $\dot{\lambda} = -a\lambda$, and $\dot{x} - ax - u = \theta N\lambda$. As before, with $x(T)$ unconstrained, this gives $\lambda(t) = \lambda(0)e^{-at}$, $\lambda(T) = \kappa/x(T)$. So

$$\begin{aligned} e^{at} \frac{d}{dt} (xe^{-at}) + \frac{1}{\lambda(0)} e^{at} &= \theta N\lambda(0) e^{-at} \\ \implies x(T)e^{-aT} - x(0) + \frac{1}{\lambda(0)} T &= \theta N\lambda(0) \frac{1}{2a} (1 - e^{-2aT}) \\ \implies (\kappa + T)u(0) - x(0) &= \frac{\theta N}{2au} (1 - e^{-2aT}). \end{aligned}$$

For $\theta = 0$ this gives the risk-neutral answer of (17.1). For $\theta > 0$ the individual should consume at a greater rate than he would do if risk-neutral, because he is optimistic that the noise will work in his favour.

Now consider what happens if $\dot{x} = ax - u$, but the lifetime T is uncertain. Let the remaining lifetime be y , with $\dot{y} = -1 + \epsilon$. The path integral is now

$$\int_0^\tau [-\log u + \lambda_2(\dot{y} + 1) + \lambda_1(\dot{x} - (ax - u)) - \frac{1}{2}\theta N\lambda_2^2] dt - \log x_\tau,$$

where τ is the first time that $y = 0$. The solution is $u = x/(\kappa + s)$, where s is the effective remaining life, which is related to x, y by

$$y^2 - s^2 = 2\theta N s^2 \left(1 - a(\kappa + s) + \log \left(\frac{\kappa + s}{x} \right) \right).$$

For a given y , this exhibits the interesting property that if x is large enough that the term in round brackets is negative, then $s > y$. In this case the individual is optimistic of a long life, and has sufficient capital to believe he can build it up, while consuming more slowly than if he were not optimistic. But if x is small then $s < y$. Having only little capital, there is not much time to build it up, and so he optimistically takes the view that his life will be short (!), and that it is best to consume what little he has at a faster rate than he would if he was certain that his remaining life were to be y .

Index

- active action, 48
- adjoint variable, 79
- asymptotically optimal, 50
- average cost, 38
- average reward, 48

- bandit process, 26, 29
- bang-bang control, 6, 76, 81
- Bellman equation, 3
- Bernoulli bandit, 29
- bin packing problem, 58
- bound, 48
- branching bandit, 36
- broom balancing, 65
- Brownian motion, 88
- Bruss's odds algorithm, 23

- calibrating bandit process, 31
- certainty equivalence, 73
- chattering, 76
- closed-loop, 3, 79
- $c\mu$ -rule, 35
- completion time, 42
- concave majorant, 24
- control theory, 1
- control variable, 2
- controllability, 63
- controllable, 63
- controlled diffusion process, 88, 89

- decomposable cost, 4
- deterministic stationary Markov policy, 20, 30
- diffusion process, 88
- discount factor, 9
- discounted programming, 11
- discounted-cost criterion, 9
- discrepancy, 94
- discrete-time, 2
- dual LP, 48
- dynamic programming equation, 3, 48

- exploration and exploitation, 18

- fair charge, 31
- feedback, 3
- finite actions, 15
- flow time, 42
- forward induction policy, 33

- gain matrix, 61
- gambling, 15
- Gittins index, 27, 30
- golf with more than one ball, 33

- Hamilton-Jacobi-Bellman equation, 74
- Hamiltonian, 79
- Hardy-Littlewood rearrangement inequality, 54
- harvesting fish, 75
- hazard rate, 35
- holding cost, 45

- index policy, 29, 30
- indexable, 49
- individual optimality, 46
- infinitesimal generator, 88
- innovations, 73
- insects as optimizers, 84
- interchange argument, 10

- job scheduling, 9

- knapsack problem, 58

- Lady's nylon stocking problem, 44
- Lagrange multipliers, 49
- linear least squares estimate, 71
- linear program, 48
- LQG model, 59

- makespan, 42
- Markov decision problem, 5
- Markov decision process, 4, 5
- Markov dynamics, 4

Markov jump process, 44
 Markov policy, 20
 monotone operator, 10
 multi-armed bandit, 29
 multi-armed bandit problem, 18, 29
 myopic policy, 18, 33

 negative programming, 11
 nonpreemptive, 35

 observability, 63
 observable, 63
 one-step look-ahead rule, 21, 33
 open-loop, 3, 79
 optimality equation, 3
 optimization of consumption, 5, 78, 96
 optimization over time, 1

 parking a rocket car, 80
 parking problem, 22
 partially observable Markov decision process, 19, 47
 passive action, 48
 perfect state observation, 4
 pharmaceutical trials, 17
 plant equation, 3
 policy, 3
 policy improvement algorithm, 40
 Pontryagin's maximum principle, 79
 positive programming, 11
 preemptive, 35
 prevailing charge, 31
 principle of optimality, 1, 2

 queueing control, 39, 45

 r -controllable, 63
 r -observable, 63
 regulation, 59
 relative value function, 39
 Riccati equation, 60, 61, 72, 75
 risk-sensitive certainly-equivalence principle, 94

 search for a moving object, 16

 secretary problem, 7, 23
 selling an asset, 12
 separable cost function, 3
 separation principle, 73
 sequential stochastic assignment problem, 53
 Shur-complement, 93
 simple family of alternative bandit processes, 26, 29

 social optimality, 46
 stability matrix, 66
 stabilizable, 66
 state variable, 3
 stochastic differential equation, 89
 stochastic scheduling, 35, 42
 stopping problem, 21
 stopping time, 27, 30
 stress, 94
 successive approximation, 15
 switching locus, 81

 target process, 34
 tax problem, 36
 threshold rule, 45
 time horizon, 2
 time to go, 5
 time-homogeneous, 5, 10
 transversality conditions, 82, 83
 turnpike theory, 87
 two-armed bandit, 27
 two-armed bandit problem, 18

 uniformization, 42, 44

 value function, 6
 value iteration, 15
 value iteration algorithm, 40
 value iteration bounds, 39

 weighted flow time, 35
 Weitzman's problem, 32
 white noise, 61
 Whittle index, 49
 Whittle index policy, 49
 Wiener process, 88