

Part I - Foundations

This material will be published by Cambridge University Press as “Markov Decision Processes and Reinforcement Learning” by Martin L. Puterman and Timothy C. Y. Chan. This pre-publication version is free to view and download for personal use only. Not for re-distribution, re-sale, or use in derivative works. ©Martin L. Puterman and Timothy C. Y. Chan, 2025.

Nature makes trees put down deep roots before having them bear fruit, and even this is done gradually^I.

Saint Vincent de Paul, French Catholic Priest, 1581-1660.

This part of the book contains two chapters:

- Chapter 2: Markov decision process fundamentals
- Chapter 3: Examples and applications

These chapters provide an important foundation for the rest of the book. As the above quote suggests, developing this foundation is critical to understanding and using the methods presented subsequently.

Chapter 2 introduces and discusses the core elements of a Markov decision process including basic components, derived quantities, and optimality criteria. Moreover, it introduces two simple Markov decision process models that provide a foundation for much of the material in this book:

1. The simple *one-period model* that highlights the fundamental trade-off between immediate and future rewards. It lays the groundwork for understanding a key feature of more complex, multi-period decision models.
2. The *two-state model* is a numerical example of a Markov decision process with two states. The two-state model recurs frequently throughout the book to illustrate key ideas and algorithms.

^Ide Paul [1995].

Familiarity with these models is crucial. They provide intuitive and concrete illustrations of the main Markov decision process and reinforcement learning concepts.

To make the model components in Chapter 2 concrete, Chapter 3 presents a broad range of illustrative applications drawn from diverse disciplines. Each explicitly identifies the key model components: states, actions, transition probabilities and rewards. As readers engage with these examples, they are encouraged to reflect on examples from their experience that can be formulated as Markov decision processes.

A strong grasp of model fundamentals is essential for successfully applying MDPs in new domains. For this reason, readers must take the time to carefully study the examples in Chapter 3, paying particular attention to how the abstract model components are defined in each context.

Chapter 2

Markov Decision Process Fundamentals

This material will be published by Cambridge University Press as “Markov Decision Processes and Reinforcement Learning” by Martin L. Puterman and Timothy C. Y. Chan. This pre-publication version is free to view and download for personal use only. Not for re-distribution, re-sale, or use in derivative works. ©Martin L. Puterman and Timothy C. Y. Chan, 2025.

Any artisan who wishes to do his job well must first sharpen his tools¹.

Confucius, Chinese philosopher, 551-479 BCE.

To effectively apply Markov decision processes and reinforcement learning, one must establish a solid grounding in the fundamental concepts and tools introduced in this chapter.

The discussion begins with a description of the core components of a Markov decision process: decision epochs, states, actions, transition probabilities, and rewards. From these elements, decision rules, policies, stochastic processes and reward processes are derived. Although not formally part of the model definition, these constructs arise naturally and inherit structure from the model’s core components. Building on them, the chapter then describes optimality criteria, which are the basis for comparing the quality of decisions and underpin the concepts of optimal policies, policy and optimal value functions, state-action value functions and the Bellman equation.

The book primarily focuses on discrete time models with finite discrete state and action spaces. It briefly discusses generalizations where appropriate such as the case of partially observable Markov decision processes (Chapter 8) in which continuous state spaces arise naturally.

¹Confucius [2003].

2.1 Basic model components

A Markov decision process model describes the fundamental elements of a recurring decision problem in which today's decision impacts future options. Decisions take place over a *planning horizon*. *Decision epochs* represent specific points of time when a decision can be made (Figure 2.1). Time is divided into *periods* or *stages*; a period begins at one decision epoch and ends at the next decision epoch.

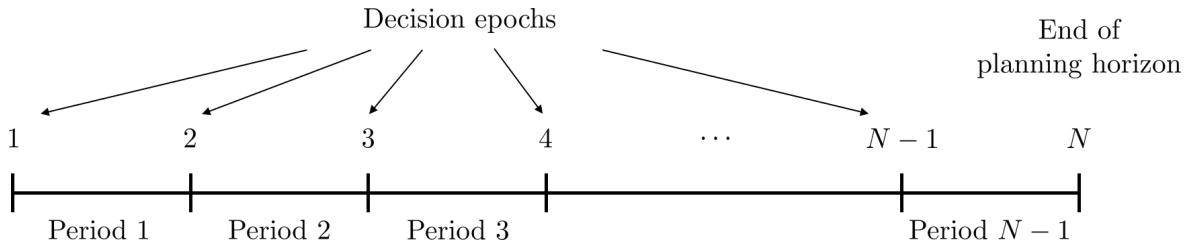


Figure 2.1: Decisions are made over a planning horizon divided into periods. The planning horizon may be finite (as shown) or infinite.

At a decision epoch, the decision maker² observes the state of the system, chooses an *action* and as a result of this choice, the system evolves to a new state according to a probability distribution that depends on the current state and the choice of action. After the system moves to the subsequent state, the decision maker receives a reward that depends on the starting state, action and possibly the subsequent state³. These steps are repeated at each subsequent decision epoch. Figure 2.2 illustrates the timing of these events between two decision epochs. Such timelines are fundamental to developing a Markov decision process model. These concepts are formalized in the following subsections.

2.1.1 Planning horizons and decision epochs

The *planning horizon* is the time interval over which decisions are made. It may be *finite* or *infinite*. A *discrete time*⁴ Markov decision process model is either:

- A *Finite Horizon Model*, in which the planning horizon is a bounded interval divided into $N - 1$ periods, or

²This book primarily uses the expression *decision maker* to refer to a possibly animate entity who makes decisions. In the computer science literature, the decision maker is often referred to as an *agent* and in the engineering literature as a *controller* reflecting a reality in which an inanimate entity in some way chooses and executes decisions.

³The reinforcement learning literature uses the acronym SARSA to refer to the sequence: *state, action, reward, state, action*.

⁴This book will not describe continuous time models.

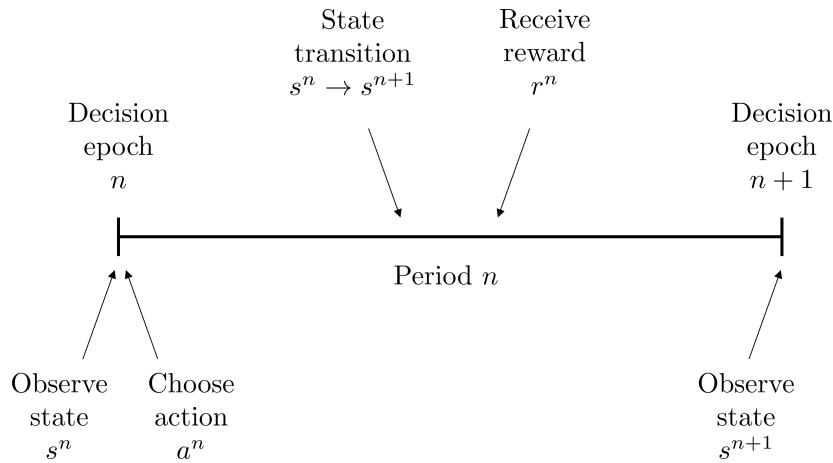


Figure 2.2: Timeline showing key events in period n . The decision maker observes state s^n at decision epoch n and chooses action a^n . The state transitions to state s^{n+1} according to the transition probabilities, the decision maker receives a reward r^n (that could depend on s^{n+1}), after which the decision maker is in the position to choose a new action at decision epoch $n + 1$ after observing state s^{n+1} .

- An *Infinite Horizon Model*, in which the planning horizon is unbounded and divided into an infinite number of periods.

Each period begins at a *decision epoch*, a point in time at which the decision maker observes the system state and chooses an action. A period ends immediately before the subsequent decision epoch. For example, each day at midnight, a retailer's inventory management system observes the stock level of all products and decides how many additional units of each to order from suppliers. In this case, the period is a day, and the decision epoch corresponds to midnight of that day.

When the planning horizon is of finite length, it is divided it into $N - 1$ periods with $N - 1$ decision epochs. No decision is made at the end of the planning horizon, epoch N , which is referred to as the *terminal epoch*. It is included to evaluate the consequences of the decision at epoch $N - 1$. Using $N - 1$ also leads to cleaner formulas. Accordingly, the book uses the convention that a finite horizon problem is referred to as an $(N - 1)$ -period problem. In a limited number of finite horizon applications such as shortest path problems, decision epochs may index the order in which decisions are made, rather than pre-defined points in time. Chapter 4 focuses on finite horizon models, while Chapters 5 – 7 focus on infinite horizon models.

Throughout the book, “realizations”, that is, states visited, actions implemented or rewards received following a decision or terminal epoch will be represented by **superscripts**. Moreover, transition probability functions and reward functions will be indexed by subscripts corresponding to the epoch to which they pertain.

In contrast, specific states and actions within the set of possible states and actions will be represented by **subscripts**^a

^aFor example a^n will denote the action chosen at decision epoch n and a_m will denote the m -th action from a set of actions.

2.1.2 States

The state of the system contains all *relevant* information available to the decision maker at a decision epoch. This means that by knowing this information and choosing an action, the transition probabilities and rewards are fully specified. Let S denote the set of all states (the *state space*) and s denote a specific *state* in S . States are necessarily *exhaustive* and *mutually exclusive*. At each decision epoch, the system must occupy a unique state in S ⁵.

Elements of S may be scalar (e.g., the inventory level of a single product) or vector-valued (e.g., queue lengths of each priority class in a multi-class queuing system). They may be ordered (e.g., the number of people in a queue) or abstract (e.g., the configuration of a Tetris screen and the shape awaiting placement in the wall).

Although our primary focus in this book is on models with S discrete and finite, there are many possible generalizations, including choosing S to be countably infinite or a subset (either bounded or unbounded) of finite-dimensional Euclidean space. Also, specifically in network or decision analysis settings, S may vary with the decision epoch, in which case one may use an epoch-specific state space S_n at epoch n . An alternative is simply to define S as the union of S_n over all decision epochs n , though this approach may unnecessarily expand the state space at each decision epoch (because some states may not be reachable at certain epochs).

The book assumes that S is independent of the decision epoch.

2.1.3 Actions

Denote by A_s the set of all actions available to the decision maker in state $s \in S$. Refer to A_s as an *action set* and each $a \in A_s$ as an *action*. Actions are mutually exclusive; the decision maker cannot choose two actions at the same time, although, as seen later, the decision maker may specify a distribution over the set of actions.

Like the state space, each A_s may be a finite set, a countably infinite set, a subset of finite-dimensional Euclidean space, or an abstract set. The action set may be inde-

⁵Unlike Schrödinger's cat, the system cannot occupy two different states at the same time.

pendent of s , as in an infinite-capacity queuing control model in which the action is to decide whether to admit or not admit a job. The book will focus on A_s being discrete and finite. Assuming the states and actions are ordered, $a_{s,j}$ represents the j -th action in the s -th state.

Note that in some states, A_s may contain a single element corresponding to “no action”. Such states are *non-actionable*. These situations occur most often when the system reaches an absorbing state or set of states, Δ , from which it cannot exit. In that case, $A_\Delta = \{\text{Do nothing}\}$. This applies to episodic models (such as controlling a robot to move from an origin to a destination or in problems in which the decision maker can stop the system at any time).

2.1.4 Transition probabilities

Let $p_n(j|s, a)$ denote the probability that the system state becomes j at decision epoch $n+1$ when the decision maker chooses action a in state s at decision epoch n . Note that several random events may occur throughout the period between decision epochs; the transition probability does not explicitly represent each one, but instead represents the net change in state. The transition probability function⁶ $p_n(j|s, a)$ has the following properties:

$$\begin{aligned} p_n(j|s, a) &\geq 0 \quad \text{for all } j \in S, a \in A_s, s \in S \\ \sum_{j \in S} p_n(j|s, a) &= 1 \quad \text{for all } a \in A_s, s \in S. \end{aligned}$$

Note that in a non-actionable state, δ , $p_n(\delta|\delta, \text{“Do nothing”}) = 1$.

When the transition probabilities do not vary over decision epochs, they are referred to as *stationary*. For infinite horizon models presented in this book, transition probabilities are assumed to be stationary, so that the subscript n is dropped.

2.1.5 Rewards

Consider two representations for a reward: $r_n(s, a, j)$ and $r_n(s, a)$. The quantity $r_n(s, a, j)$ denotes the reward received in period n when the decision maker chooses action a in state s and the system transitions to state j at decision epoch $n+1$. The latter quantity, $r_n(s, a)$, has a similar interpretation, but is independent of the subsequent state. The choice of reward function depends on the application – usually one of these two forms better describes a specific context. The examples in Chapter 3 will illustrate both.

Assume $r_n(s, a, j)$ and $r_n(s, a)$ are scalar and real-valued. Generalizations include vector-valued or abstract rewards. For example, as the result of an action in a fantasy

⁶When S represents an uncountable subset of finite-dimensional Euclidean space, the transition probability may be represented by a probability density function.

game, the decision maker may receive a sword, a shield and a vial of magic potion. Instead of keeping track of this bundle of goods in the reward function, the decision maker might assign a numerical value to each item and be indifferent between collections of items with the same total value.

These quantities are *rewards* because the book formulates the problem of a decision maker seeking to *maximize* rewards. *Costs* can be represented as negative rewards to allow for a decision maker who seeks to minimize costs. Consequently, minimizing costs corresponds to maximizing rewards. To avoid awkward minus signs in expressions, on some occasions when minimizing costs, rewards will be replaced by $c_n(s, a)$ to represent the cost of being state s and choosing action a or $c_n(s, a, j)$ to include the possibility that the cost depends on the state at the next decision epoch j .

When using optimality criteria based on expected rewards⁷, the quantity $r_n(s, a)$ may also represent the *expected* reward in period n where the expectation is taken over the possible states at decision epoch $n + 1$:

$$r_n(s, a) = \sum_{j \in S} r_n(s, a, j) p_n(j|s, a). \quad (2.1)$$

For discrete time models, the model formulation does not account for how the reward is accumulated throughout a period between two decision epochs. For example, in an inventory control model with weekly decision epochs, the inventory levels may change during the week resulting in varying holding costs between decision epochs. A discrete time formulation accumulates all of these costs and summarizes them in the reward function for that one period.

In finite horizon models, a *terminal reward* or *scrap value* $r_N(s)$ is used to represent the consequence of ending the planning horizon in state s . In infinite horizon models, the terminal reward is omitted. Furthermore, the subscript n is deleted from the reward function in the infinite horizon setting, since the focus in that case is on *stationary* rewards.

Rewards are an intrinsic part of the model that arise naturally from the application. However, in reinforcement learning models, the modeler may need to construct an *artificial* reward function that conforms with the objectives of the task.

2.1.6 Markov decision processes and decision trees

A decision tree provides a visual display of the basic model components. In Figure 2.3 square boxes represent states, arcs from boxes to circles represent possible actions in each state, arcs from circles to subsequent states denote transitions that occur according to a transition probability and result in a reward.

⁷In some applications, particularly in economic contexts, a decision maker seeks to maximize expected *utility* or some other risk sensitive criterion.

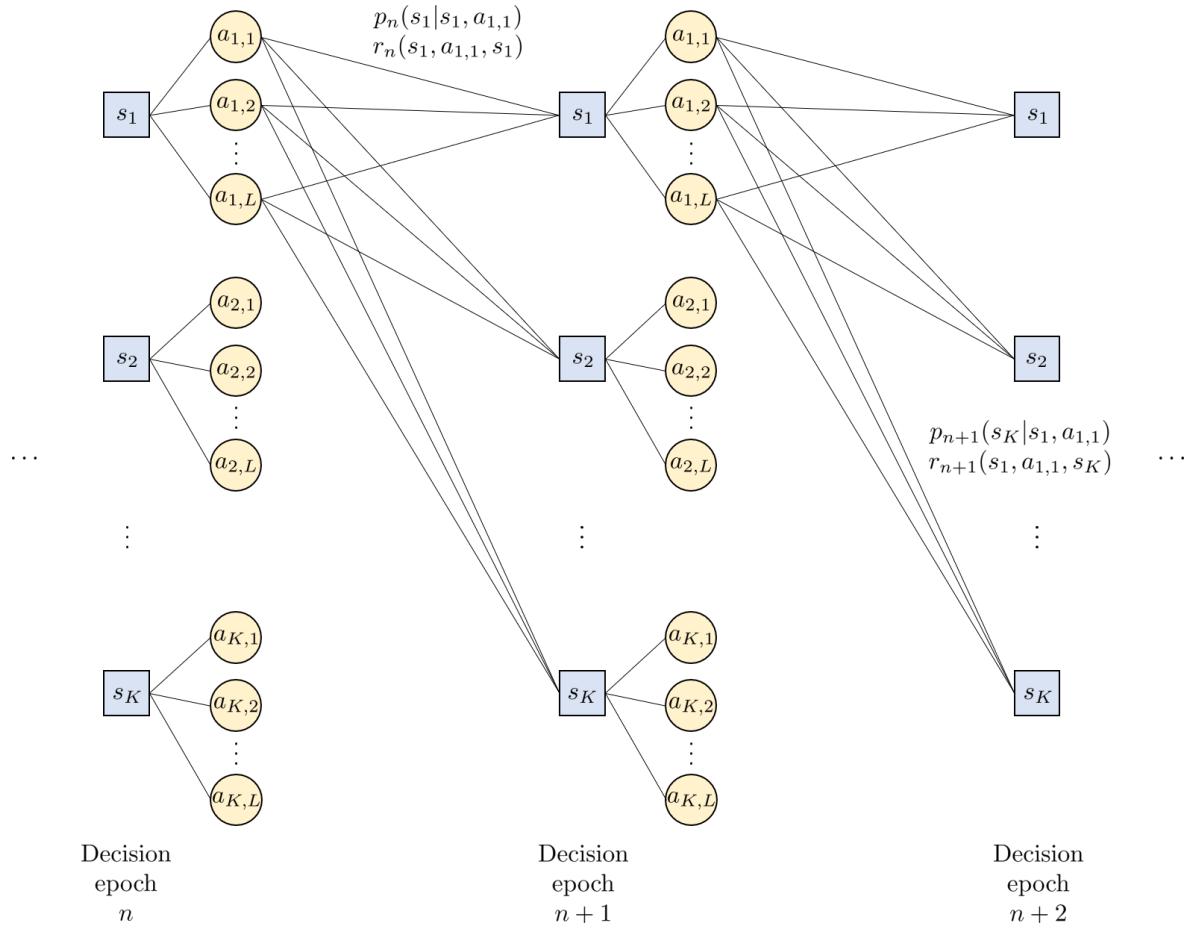


Figure 2.3: Graphical representation of a Markov decision process by a decision tree.

It should be evident from the decision tree representation that, in general, there is an exponential explosion in the number of possible trajectories through the tree as the length of the planning horizon is increased.

2.2 Derived objects

A Markov decision process is fully specified⁸ by the five model components described in the previous section:

- the planning horizon N (which may be finite or infinite),
- the set of states S ,

⁸Contrary to formulations in some disciplines, the discount rate is not specified as a model component. It is more appropriately associated with the optimality criterion.

- the sets of actions A_s for each $s \in S$,
- the transition probabilities $p_n(j|s, a)$, and
- the rewards $r_n(s, a, j)$ or $r_n(s, a)$.

In this section, decision rules, policies, derived stochastic processes and reward processes are described. They are not part of the basic formulation, but are fundamental concepts derived from the basic model components.

2.2.1 Decision rules

A *decision rule* describes both the information and mechanism a decision maker uses to select an action in a given state at a single, specific decision epoch. Decision rules can be classified on the basis of two independent dimensions:

Information: Markovian or history-dependent

Mechanism: Deterministic or randomized

Information

The Markovian vs. history-dependent dichotomy describes the *information* used by the decision maker when choosing an action. A *Markovian* decision rule uses only the state at the current decision epoch to select actions, while a *history-dependent* decision rule considers some or all of the previous states and actions up to and including the current state when choosing an action. That is, at decision epoch n , a Markovian decision rule is a function of s^n and a history-dependent decision rule is a function of a *trajectory* $(s^1, a^1, s^2, \dots, a^{n-1}, s^n)$ ⁹. Thus, a Markovian decision rule is a special case of a history-dependent decision rule in which the history is summarized in a single state. Write

$$H_n = \{h^n = (s^1, a^1, s^2, \dots, a^{n-1}, s^n) \mid s^1 \in S, a^1 \in A_{s^1}, s^2 \in S, \dots, a^{n-1} \in A_{s^{n-1}}, s^n \in S\} \quad (2.2)$$

as the set of all histories, h^n , leading up to decision epoch n . Note that $H_1 = S$. While the past sequence of rewards could also be explicitly included in the history, its inclusion is redundant in the Markov decision process model, since the history of states and actions alone allows one to reconstruct the past rewards through the reward functions¹⁰. Note that $h^1 = s^1$ and for $n = 2, 3, \dots, N$,

$$h^n = (h^{n-1}, a^{n-1}, s^n). \quad (2.3)$$

This recursion is used explicitly in Section 4.1.3.

⁹Recall the convention that superscripts refer to the state and/or action chosen at decision epoch n . Also, history-dependent decision rules need not consider the entire history, but a subset of past actions and states.

¹⁰In model-free reinforcement learning, a functional form for $r_n(s, a, j)$ is not known or specified. In this case the history may be augmented to include the reward.

Mechanism

The deterministic vs. randomized dichotomy describes the *mechanism* used to select an action at a given decision epoch. A *deterministic* decision rule selects an action with certainty, while a *randomized* decision rule selects an action according to a specified probability distribution. A deterministic decision rule is a special case of a randomized decision rule corresponding to a degenerate probability distribution¹¹.

Types of decision rules

The two dimensions described above lead to four classes of decision rules.

- A *Markovian deterministic* (MD) decision rule, d_n , is a function from states to actions. More formally, $d_n(s^n) = a^n$ denotes the decision rule that at decision epoch n chooses action $a^n \in A_{s^n}$ when the system is in state $s^n \in S$.
- A *history-dependent deterministic* (HD) decision rule, d_n , is a function from the set of histories to actions. More formally, $d_n(h^n) = a^n$ denotes the decision rule that chooses action $a^n \in A_{s^n}$ when the system history is $h^n \in H_n$ and s^n is the state at decision epoch n . The subtlety here is that although the decision rule's action choice may vary with history, it can only choose actions from the set A_{s^n} at epoch n , which depends only on the state at decision epoch n . This means that given two different histories h^n and \hat{h}^n that both arrive at state s^n , $d_n(h^n)$ and $d_n(\hat{h}^n)$ may choose different actions, but each action will be chosen from the same action set A_{s^n} .
- A *Markovian randomized* (MR) decision rule, d_n , is a mapping from the set of states to the set of probability distributions over the action set. More formally, $w_{d_n}^n(a^n|s^n)$ denotes the probability that decision rule d_n chooses action a^n in state s^n .
- A *history-dependent randomized* (HR) decision rule, d_n , is a mapping from the set of histories to the set of probability distributions over the action set. More formally, $w_{d_n}^n(a^n|h^n)$ denotes the probability that decision rule d_n chooses action a^n given history h^n .

These classes of decision rules are denoted as D^{MD} , D^{HD} , D^{MR} and D^{HR} , respectively.

It is important to note that because of the ability to construct history-dependent decision rules, the Markov decision process formulation gives rise to a system that might not evolve in a Markovian fashion. The expression “Markov decision process” refers to the fact that rewards and transition probabilities depend on the past only through the state and action at the present decision epoch. It does not, however, mean

¹¹A degenerate probability distribution places all of the probability on one action.

that every stochastic process generated by this model is a Markov chain. See Section 2.2.3 below for more details on this point.

Subsequent chapters will establish the optimality of Markovian deterministic decision rules so it will be unnecessary to consider randomized decision rules when seeking optimal policies. However, randomized decision rules are fundamental in linear programming formulations of the infinite-horizon Markov decision process model (Chapters 5–7), exploration-based simulation algorithms (Chapter 10) and policy-based reinforcement learning algorithms (Chapter 11).

2.2.2 Policies

A *policy* π , sometimes referred to as a *contingency plan* or *strategy*, is a sequence of decision rules, one for each decision epoch. In finite horizon models, write $\pi = (d_1, d_2, \dots, d_{N-1})$. In infinite horizon models, $\pi = (d_1, d_2, \dots)$. In the classical Markov decision process setting, once an optimality criterion has been specified, the decision maker's goal is to find an optimal policy – one that maximizes or minimizes the designated criterion.

The four classes of decision rules defined in Section 2.2.1 result in four classes of policies, Π^{MD} , Π^{HD} , Π^{MR} and Π^{HR} . In addition, stationary deterministic (SD) and stationary randomized (SR) policies, denoted by the classes Π^{SD} and Π^{SR} , respectively, refer to sets of *stationary policies*, that is policies¹² that choose the same decision rule at every decision epoch. In infinite horizon models, stationary policies that use the decision rule d at every decision epoch will be denoted by d^∞ , that is,

$$d^\infty := (d, d, \dots). \quad (2.4)$$

Some comments about policies follow:

- The class Π^{HR} denotes the most general class of policies; all policies considered in the book are in this class. This level of generality is required to define optimality, but often not required to find optimal policies. For example, in finite horizon models it will be sufficient to restrict one's search for an optimal policy to the class of Markovian deterministic policies (see Section 2.4).
- Policies in Π^{HR} generally cannot be implemented in long finite-horizon or infinite horizon models because storage requirements increase exponentially with respect to the planning horizon length. Fortunately, in finite horizon models, there exist optimal policies in this class that are Markovian and deterministic, and in infinite horizon models, there exist optimal policies that are stationary and deterministic.
- Stationary policies are most relevant to infinite horizon models. Proofs and methods in Chapters 5–6 will use the result that there exist stationary deterministic

¹²In most cases, a stationary policy is necessarily Markovian.

policies that are optimal in the class of history-dependent randomized policies under several different optimality criteria.

- In finite horizon models, a Markovian deterministic policy may be characterized by a *lookup table*. A lookup table is an array in which rows correspond to states, columns correspond to decision epochs and an entry indicates which action the policy selects in that state at that decision epoch. Although a useful conceptual framework, a lookup table may be an impractical method of encoding a policy in models with a large number of states. Chapters 9 and 11 replace lookup tables by function approximation. Models in which lookup tables apply are referred to as *tabular models* in the reinforcement learning literature.
- Figure 2.4 summarizes the relationship between policy classes¹³.

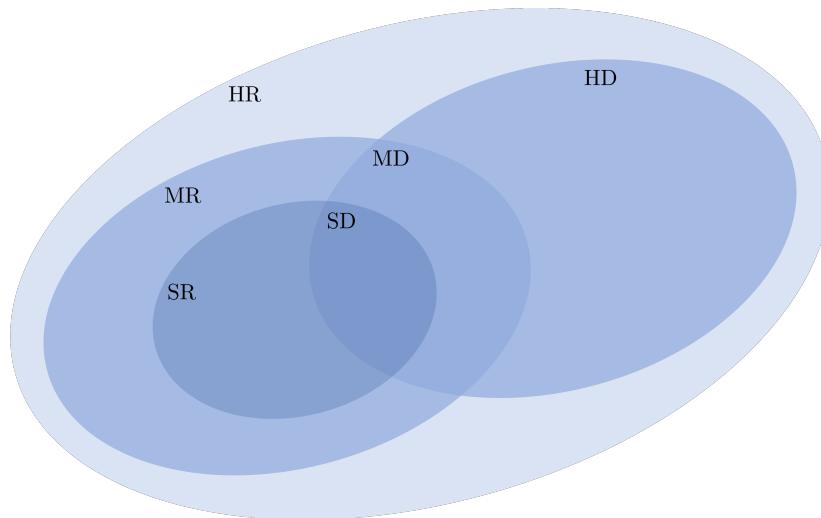


Figure 2.4: Relationships between policy classes. Labels near the edge of an oval refer to the entire oval, whereas labels inside the intersection of two ovals denote the entire intersection between them. For example, SD is the intersection of SR and HD. MD is the intersection of MR and HD, which includes SD.

2.2.3 Derived stochastic processes

Once a policy is chosen and a probability distribution over the starting state of the process is specified, the probabilistic evolution of a Markov decision process is completely

¹³Although not represented in the figure, there may exist policies that are stationary and history-dependent (i.e., non-Markovian). An example is a policy that makes decisions based on the current state as well as the starting state. However, these examples represent “edge cases” and are uncommon. Thus, in this book, only stationary policies are considered within the Markovian class.

determined. Let the random variable X_1 denote the initial state of the Markov decision process and let $\rho(s) := P[X_1 = s]$ denote the initial state probability distribution, which by assumption is independent of policy choice. When, as in most applications, the system starts in a specific state at decision epoch 1, say s^1 , write $\rho(s) = 1$ for $s = s^1$ and $\rho(s) = 0$ for $s \neq s^1$.

Let π denote a policy for either a finite or infinite horizon problem and let the random variable Y_1 denote the action chosen by d_1 at decision epoch 1. There are two potential sources of variability that affect Y_1 :

1. Variability in X_1 , which occurs when X_1 is random, and
2. Variability in action choice in X_1 , which occurs if d_1 is a randomized decision rule.

Once Y_1 is chosen, the next state, X_2 , is random with a distribution determined by the transition probabilities. This is true even if X_1 and Y_1 are not random, that is when the system starts in a specific state s^1 and d_1 is deterministic. Consequently, the second action Y_2 will be random as well, and so on. Thus, a policy and initial state distribution generate the stochastic process

$$(X_1, Y_1, X_2, Y_2, \dots, X_{N-1}, Y_{N-1}, X_N) \quad (2.5)$$

in a finite horizon model and

$$(X_1, Y_1, X_2, Y_2, \dots) \quad (2.6)$$

in an infinite horizon model. In these expressions, X_n represents the state of the stochastic process and Y_n denotes the action chosen by policy π at decision epoch n .

Some comments about policies with respect to derived stochastic processes follow:

1. When $\pi \in \Pi^{\text{MR}}$, the stochastic process is a discrete time Markov chain. See Exercise 3.
2. When π is history-dependent, the stochastic process need not be a Markov chain. That is, at each decision epoch, given a state and action, the transition probabilities may depend on all or some of the past. See Exercise 3.
3. A policy $\pi = (d_1, d_2, \dots)$ induces a probability distribution p^π on realizations of the stochastic process $(X_1, Y_1, X_2, Y_2, \dots)$. These realizations are referred to as *trajectories*. To determine this distribution requires enumerating realizations and multiplying the transition probabilities and action randomization distributions corresponding to each realization as follows:
 - (a) For $\pi \in \Pi^{\text{HR}}$,

$$p^\pi(s^1, a^1, s^2, a^2, s^3, \dots)$$

$$= \rho(s^1)w_{d_1}^1(a^1|s^1)p_1(s^2|s^1, a^1)w_{d_2}^2(a^2|h^2)p_2(s^3|s^2, a^2)w_{d_3}^3(a^3|h^3)\dots \quad (2.7)$$

where $h^1 = s^1$, $h^2 = (s^1, a^1, s^2)$, and so on. Recall that $w_{d_n}^n(a|h)$, defined in Section 2.2.1, denotes the probability that decision rule d_n chooses action a in state s at epoch n given history h . That is $w_{d_n}^n = P^\pi[Y_n = a]$ for $\pi = (d_1, d_2, \dots)$.

Note that the entire history first enters into this expression through the randomization distribution at decision epoch 2 and subsequently, but not through the transition probabilities which only depend on the state and action. If a deterministic policy was used, the history would also only enter through the decision rule as in (2.7), but the decision rule would explicitly choose the action in each transition probability.

(b) For $\pi \in \Pi^{\text{MR}}$,

$$\begin{aligned} p^\pi(s^1, a^1, s^2, a^2, s^3, \dots) \\ = \rho(s^1)w_{d_1}^1(a^1|s^1)p_1(s^2|s^1, a^1)w_{d_2}^2(a^2|s^2)p_2(s^3|s^2, a^2)w_{d_3}^3(a^3|s^3)\dots \end{aligned} \quad (2.8)$$

This expression is similar to (2.7), except that histories are replaced with states.

(c) For $\pi \in \Pi^{\text{MD}}$

$$p^\pi(s^1, a^1, s^2, a^2, s^3, \dots) = \rho(s^1)p_1(s^2|s^1, a^1)p_2(s^3|s^2, a^2)\dots \quad (2.9)$$

since $d_1(s^1) = a^1$, $d_2(s^2) = a^2$, and so on. In this case, transition probabilities need only be multiplied with each other (and the initial state distribution) to find the probability distribution of the stochastic process generated by π .

2.2.4 Reward processes

Section 2.2.3 showed that a policy π generates a stochastic process of states and actions represented by (2.5) in finite horizon models and (2.6) for infinite horizon models with distribution $p^\pi(\cdot)$ defined in (2.7) in the most general case. These processes generate corresponding stochastic processes of rewards

$$(r_1(X_1, Y_1, X_2), r_2(X_2, Y_2, X_3)\dots, r_{N-1}(X_{N-1}, Y_{N-1}, X_N), r_N(X_N)) \quad (2.10)$$

in finite horizon models, and

$$(r_1(X_1, Y_1, X_2), r_2(X_2, Y_2, X_3), \dots) \quad (2.11)$$

in infinite horizon models.

In greater generality, let these stochastic processes be represented by either

$$\mathcal{R}_N := (R_1, R_2, \dots, R_N) \quad (2.12)$$

or

$$\mathcal{R}_\infty := (R_1, R_2, \dots) \quad (2.13)$$

where the random variable R_n represents the reward received in period n .

Letting r^n denote a realization of R_n , the expressions

$$P[\mathcal{R}_N = (r^1, r^2, \dots, r^n)] \quad \text{and} \quad P[\mathcal{R}_\infty = (r^1, r^2, \dots)]$$

equal the probability of the given finite or infinite sequence of rewards. When reward processes arise in a Markov decision process, the policy-dependent probability $p^\pi(\cdot)$ replaces the generic probability $P[\cdot]$.

Simulation or observing repeated experiments in the real world will generate many different sequences of rewards. The probability distribution of \mathcal{R}_N and \mathcal{R}_∞ describes the likelihood of observing each sequence. In this book, there will be little occasion to explicitly compute such probability distributions, but they will provide the basis for interpreting several expressions. Most calculations will avoid this onerous task by evaluating distributions and expectations sequentially. The example in Section 4.1.2 shows how to derive the distribution of \mathcal{R}_N in an Markov decision process by determining which states and actions map into the same reward sequence.

In the discussion of values in the next section, it is sufficient to work directly with the reward process described above. However, in the simulation methods in Chapters 10 and 11, realizations (trajectories) that include states, actions and rewards of the form^[14]

$$(X_1, Y_1, R_1, X_2, Y_2, \dots, R_N) \quad \text{and} \quad (X_1, Y_1, R_1, X_2, Y_2, \dots)$$

will be considered. Thus a realization of the process will be written as either

$$(s^1, a^1, r^1, s^2, \dots, r^N) \quad \text{or} \quad (s^1, a^1, r^1, s^2, \dots).$$

Following convention, a stochastic process together with a reward function will be referred to as a *reward process*. When the stochastic process is a Markov chain, it is a *Markov reward process*. Markov reward processes most often arise as the sequence of rewards of a Markovian policy in a Markov decision process.

2.2.5 Assigning a value to a reward process

This section shows how reward processes lead to decision making criteria. As noted in Section 2.1.5, assume that reward functions $r_n(s, a, j)$ are real-valued so that each R_n , as defined implicitly through (2.12) or (2.13), is real-valued. Thus, the random reward sequence $\mathcal{R}_N = (R_1, R_2, \dots, R_N)$ takes values in \mathbb{R}^N ^[15] and $\mathcal{R}_\infty = (R_1, R_2, \dots)$ takes values in \mathbb{R}^∞ .

¹⁴Technically speaking when $r_n(\cdot, \cdot, \cdot)$ is a function of the current state, the current action and the next state, we need to observe X_{n+1} before we can evaluate the reward. So technically, at least, we should write the stochastic process $(X_1, Y_1, X_2, R_2, Y_2, \dots, X_N, R_N)$.

¹⁵Denotes N -dimensional Euclidean space.

Expected utility

To assign a value to a sequence of rewards, a decision maker may use *expected utility*. Let $u(\cdot)$ denote a real-valued function on \mathbb{R}^N or \mathbb{R}^∞ , and let $E[\cdot]$ denote expectation with respect to the probability distribution of \mathcal{R}_N or \mathcal{R}_∞ . The expected utility of these reward sequences equals $E[u(R_1, R_2, \dots, R_N)]$ and $E[u(R_1, R_2, \dots)]$, respectively.

Often utility functions are additive, so

$$u(R_1, R_2, \dots, R_N) = \sum_{n=1}^N u_n(R_n) \quad (2.14)$$

for N finite or infinite, where $u_n(\cdot)$ is a real-valued function for each n . When N is finite, the expected utility becomes

$$E[u(\mathcal{R}_N)] = E[u(R_1, R_2, \dots, R_N)] = E \left[\sum_{n=1}^N u_n(R_n) \right]. \quad (2.15)$$

When N is infinite, define¹⁶

$$E[u(\mathcal{R}_\infty)] = E[u(R_1, R_2, \dots)] := \lim_{N \rightarrow \infty} E \left[\sum_{n=1}^N u_n(R_n) \right]. \quad (2.16)$$

Interpreting utilities

Utilities $u_n(\cdot)$ encode a decision maker's attitude towards both risk (variability) and timing of rewards. The most common choice for $u_n(\cdot)$ is

$$u_n(r) = \lambda^{n-1} r$$

where $0 \leq \lambda \leq 1$ so that

$$E[u_n(R_n)] = \lambda^{n-1} E[R_n].$$

Refer to the quantity λ as the *discount factor* or *discount rate*. When $\lambda = 1$ the decision maker is indifferent to the timing of rewards. When $\lambda < 1$ the decision maker prefers receiving rewards sooner than later. For example, if $\lambda = 0.95$, a decision maker would be indifferent between receiving one unit of utility next period and 0.95 units now.

When using $E[R_n]$ to make decisions, a decision maker is said to be *risk-neutral*. For example, a risk-neutral decision maker would be indifferent between the following two gambles since they have the same expected value:

- Gamble A: win \$100 with probability 1, or

¹⁶In (2.16), it would be preferable to take the limit inside the expectation but that limit need not exist. An example below explores this point further

- Gamble B: win \$0 with probability 0.5 and \$200 with probability 0.5.

A *risk-seeking* decision maker would prefer Gamble B and a *risk-averse* decision maker would prefer Gamble A. Convex increasing utility functions such as r^2 , correspond to risk-seeking decision making while concave increasing utility functions such as \sqrt{r} correspond to risk-averse decision making¹⁷.

As an example, consider coaching decisions in sport. If a team is behind near the end of the game, its coach may be risk-seeking in an attempt to catch up. Alternatively, if a team is ahead, its coach may choose to make more conservative decisions, based on a risk-averse utility function, in order to hold on to the lead.

The Markov decision process models considered in this book will focus on decision making by risk-neutral decision makers, so utility is replaced by rewards or discounted rewards directly. However, it is important to be aware that other utility functions apply, and optimal policies with respect to one utility function will not necessarily be optimal for a different utility function.

2.3 Optimality criteria: Transforming a Markov decision process into a Markov decision problem

Up to this point, the Markov decision process has been formulated without considering the decision maker's preferences for reward sequences generated by different policies. The following sections define and comment on the following three optimality criteria that are commonly used to evaluate and compare policies:

- expected total reward,
- expected discounted reward, and
- long-run average reward.

The expected total reward criterion applies to both finite and infinite horizon models, while the latter two are most often used in infinite horizon models. A policy is said to be *optimal* when it maximizes the appropriate criterion over the set of all history-dependent randomized policies. To be precise, the expression *Markov decision problem* will correspond to a Markov decision process *together* with an optimality criterion¹⁸. However, as with most of the published literature, this book will use the general phrase *Markov decision process* to refer to both cases, whether an optimality criterion is included or not.

¹⁷Which would better reflect your attitude towards gambles? Why?

¹⁸We emphasize that a Markov decision process is defined independent of the optimality criterion.

2.3.1 Expected total reward

This section defines the expected total reward of a reward sequence and policy. Technical considerations lead us to distinguish finite horizon from infinite horizon models.

Finite horizon

Define the *expected total reward* of a finite sequence of random rewards (R_1, R_2, \dots, R_N) by

$$E \left[\sum_{n=1}^N R_n \right] \quad (2.17)$$

where the expectation is over the distribution of reward sequences. Since many reward sequences might map into the same total reward, in theory, when R_n is discrete, this quantity can be computed by first enumerating all possible values for the sum and then accumulating the probabilities of the trajectories with the same sum.

As noted in Section 2.2.3, a policy π generates a random sequence

$$(X_1, Y_1, X_2, Y_2, \dots, X_{N-1}, Y_{N-1}, X_N)$$

of states and actions and consequently generates a random sequence of rewards by setting $R_n = r_n(X_n, Y_n, X_{n+1})$ or $R_n = r_n(X_n, Y_n)$ for $n = 1, 2, \dots, N - 1$ and $R_N = r_N(X_N)$.

Thus, the expected total reward of a policy $\pi \in \Pi^{\text{HR}}$ can be defined as follows.

Definition 2.1. For all $s \in S$, the *finite horizon expected total reward* of policy π is

$$v^\pi(s) := E^\pi \left[\sum_{n=1}^N R_n \middle| X_1 = s \right]. \quad (2.18)$$

In equation (2.18), the expectation is with respect to the probability distribution of random rewards that result from different realizations of the stochastic process generated by π with $X_1 = s$ (Section 2.2.4).

Refer to the expected total reward of policy π , $v^\pi(s)$, as its *value*. Implicit in this notation for finite horizon models is that $v^\pi(s)$ gives the value for a model with $N - 1$ decision epochs and terminal epoch N ¹⁹. The function $v^\pi(\cdot)$ will be called a *policy value function*. Note that $v^\pi(s)$ will most often be seen written as either

$$v^\pi(s) = E^\pi \left[\sum_{n=1}^{N-1} r_n(X_n, Y_n, X_{n+1}) + r_N(X_N) \middle| X_1 = s \right] \quad (2.19)$$

¹⁹Some authors use the notation $v_N^\pi(s)$ where N denotes the planning horizon.

or

$$v^\pi(s) = E^\pi \left[\sum_{n=1}^{N-1} r_n(X_n, Y_n) + r_N(X_n) \middle| X_1 = s \right] \quad (2.20)$$

depending on the form of r_n . To simplify notation, $E_s^\pi[\cdot]$ will often be used instead of $E^\pi[\cdot|X_1 = s]$.

Simulating the expected total reward

The following algorithm describes how to use simulation to estimate $v^\pi(s)$. It generates a single replicate for a single starting state s^1 . Since the decision maker often needs $v^\pi(s)$ for all $s \in S$, this algorithm would be repeated for all states. *Monte Carlo estimates* (Chapter 10) are obtained averaging $v^\pi(s)$ over many replicates.

Algorithm 2.1. Simulating a policy value function for a deterministic policy based on a single replicate for a finite horizon model.

1. **Initialize:**

- (a) Specify $\pi = (d_1, d_2, \dots, d_{N-1})$.
- (b) Specify $s^1 \in S$.
- (c) $n \leftarrow 1$ and $v \leftarrow 0$.

2. **Iterate:** While $n \leq N - 1$:

- (a) $a^n \leftarrow d(s^n)$.
- (b) Sample s^{n+1} from $p_n(\cdot|s^n, a^n)$.
- (c) $r^n \leftarrow r_n(s^n, a^n, s^{n+1})$.
- (d) $v \leftarrow v + r^n$.
- (e) $n \leftarrow n + 1$.

3. $v \leftarrow v + r_N(s^N)$

4. **Terminate:** Return v .

Some comments about simulating the expected total reward follow:

1. Although you might not intend to simulate this process, this algorithm describes how the process would evolve in an implementation.
2. The functions $p_n(\cdot|\cdot, \cdot)$ and $r_n(\cdot, \cdot, \cdot)$ are explicitly used in the above procedure. When this is done, the approach is said to be *model-based*. Alternatively if s^{n+1} in step 2(b) and r^n in step 2(c) were outputs of a simulation or real-time process,

the approach is said to be *model-free*. This distinction will be most significant in Chapters 10 and 11.

3. If one wished to instead simulate a randomized policy, step 2(a) would be replaced by “Sample a^n from $w_{d_n}^n(\cdot|s^n)$ ”.
4. Obvious modifications would be required to evaluate a history-dependent policy. This would be impractical if N is large and the decision rules depended on the whole past.

Chapter 4 will develop a straightforward approach to compute $v^\pi(s)$ numerically or analytically for any $s \in S$ and any π . However, in problems with large state spaces, simulation and approximation may be preferable.

Infinite horizon models

The expected total reward (value) of a policy π in an infinite horizon model can be defined as follows.

Definition 2.2. For all $s \in S$, the *infinite horizon expected total reward* of policy π is

$$v^\pi(s) := \lim_{N \rightarrow \infty} E_s^\pi \left[\sum_{n=1}^N R_n \right]. \quad (2.21)$$

Note that the expected total reward of a policy is defined as the limit of the expected value of its finite sums. This is because

$$\lim_{n \rightarrow \infty} \sum_{n=1}^N R_n$$

need not exist so that its expectation may be undefined. Equation (2.21) can be written in terms of the reward function directly as

$$v^\pi(s) = \lim_{N \rightarrow \infty} E_s^\pi \left[\sum_{n=1}^N r_n(X_n, Y_n, X_{n+1}) \right] \text{ or } v^\pi(s) = \lim_{N \rightarrow \infty} E_s^\pi \left[\sum_{n=1}^N r_n(X_n, Y_n) \right]. \quad (2.22)$$

In contrast to the finite horizon setting, it is now of concern whether the limits in (2.22) exist. Some possible limiting behaviors for this (or any) sequence include:

Convergence: when $E[R_n]$ decreases sufficiently quickly or becomes zero eventually,

Divergence: when $E[R_n]$ remains sufficiently positive or negative,

Oscillation: when $E[R_n]$ alternates between positive and negative values and does not die out.

When using the expected total reward criterion, the infinite horizon Markov decision process literature has focused on models in which the limit in (2.21) exists. At first glance, it may not be apparent how these expectations can be finite when all rewards are positive. The reason is that most examples of this kind are either *episodic models* or *optimal stopping problems*, namely, models that terminate in *reward-free* absorbing states²⁰. Such models include policies that ensure the expected time to reach one of the absorbing states is finite so that effectively they behave like finite horizon problems but with policy-specific variable horizons. Often the random horizon is referred to as the *effective horizon*.

In such problems, the decision maker attempts to delay reaching an absorbing state as long as possible when rewards are mostly positive. When rewards are mostly negative, the decision maker attempts to reach an absorbing state as quickly as possible. These models include *transient* models and *stochastic shortest path* models. See Sections 3.5, 3.2, and 3.8 for concrete examples of such problems.

An illustrative example*

The following example justifies the above definition of the expected total reward in infinite horizon models.

Example 2.1. Consider the Markov reward process depicted in Figure 2.5. It contains three states $S = \{s_1, s_2, s_3\}$. Rewards and transition probabilities are represented in brackets below the arcs. The first number in the bracket represents the reward received if that transition takes place and the second number represents the probability of that transition. Thus, from state s_1 the system transitions to state s_2 or s_3 with probability 0.5 and generates a reward of 0. From state s_2 the Markov chain transitions to state s_3 with certainty and generates a reward of 1 and from state s_3 the system transitions to state s_2 with certainty and generates a reward of -1 . No other transitions are possible.

Next consider the calculation of $v(s_1)$. There are two possible realizations of the Markov chain describing the state at epoch n , namely

$$(s_1, s_2, s_3, s_2, s_3, \dots) \quad \text{or} \quad (s_1, s_3, s_2, s_3, s_2, \dots).$$

Each occurs with probability 0.5. The corresponding Markov reward process generates sequences

$$\xi_1 := (0, -1, 1, -1, 1, \dots) \quad \text{and} \quad \xi_2 := (0, 1, -1, 1, -1, \dots).$$

²⁰Appendix B defines basic Markov chain concepts.

each occurring with probability 0.5. Thus

$$P[\mathcal{R}_\infty = \xi_1 | X_1 = s_1] = P[\mathcal{R}_\infty = \xi_2 | X_1 = s_1] = 0.5.$$

The sequence of **total rewards** corresponding to each realization are

$$\sigma_1 := (0, -1, 0, -1, 0, \dots) \quad \text{and} \quad \sigma_2 := (0, 1, 0, 1, 0, \dots).$$

For finite N even, the total reward

$$P \left[\sum_{n=1}^N R_n = -1 \mid X_1 = s_1 \right] = P \left[\sum_{n=1}^N R_n = 1 \mid X_1 = s_1 \right] = 0.5,$$

while for N odd,

$$P \left[\sum_{n=1}^N R_n = 0 \mid X_1 = s_1 \right] = 1.$$

Next let $N \rightarrow \infty$. Since both possible total reward series oscillate, neither σ_1 or σ_2 has a limit. Hence the quantity $\lim_{N \rightarrow \infty} \sum_{n=1}^N R_n$ does not exist, so that the expression

$$E \left[\lim_{N \rightarrow \infty} \sum_{n=1}^N R_n \mid X_1 = s_1 \right]$$

is not well defined. However for any N ,

$$E \left[\sum_{n=1}^N R_n \mid X_1 = s_1 \right] = 0,$$

so

$$\lim_{N \rightarrow \infty} E \left[\sum_{n=1}^N R_n \mid X_1 = s_1 \right] = 0$$

so that the limit in (2.21) exists.

This is the reason the expected total reward in an infinite horizon model is defined as the “limit of expected partial sums” as opposed to the “expectation of the limit of partial sums”.

Note the limit can be passed inside the expectation by considering the *lim inf* and *lim sup* of the sequence of partial sums. These quantities always exist and are well defined. Thus, in this example

$$P \left[\limsup_{N \rightarrow \infty} \sum_{n=1}^N R_n = 0 \mid X_1 = s_1 \right] = P \left[\limsup_{N \rightarrow \infty} \sum_{n=1}^N R_n = 1 \mid X_1 = s_1 \right] = 0.5$$

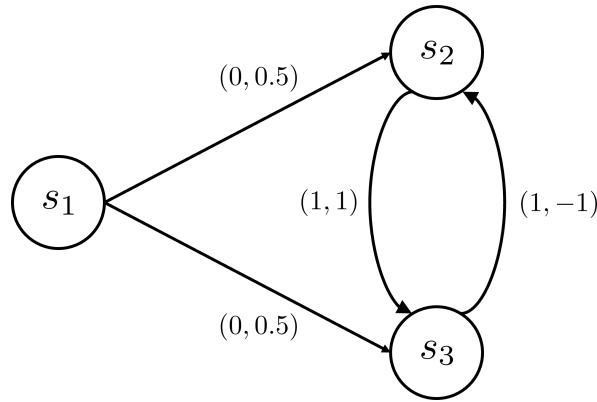


Figure 2.5: Markov reward process analyzed in Example 2.1. The labels (r, p) refer to the reward and transition probabilities, respectively.

so that

$$E \left[\limsup_{N \rightarrow \infty} \sum_{n=1}^N R_n \mid X_1 = s_1 \right] = 0.5.$$

Moreover

$$E \left[\liminf_{N \rightarrow \infty} \sum_{n=1}^N R_n \mid X_1 = s_1 \right] = -0.5.$$

Optimal policies

Definition 2.3. A policy π^* is *optimal* if

$$v^{\pi^*}(s) \geq v^\pi(s) \quad (2.23)$$

for all $s \in S$ and $\pi \in \Pi^{\text{HR}}$.

Chapter 4 analyzes finite horizon Markov decision process models and Chapter 6 considers infinite horizon models under this optimality criterion.

2.3.2 Expected discounted reward

In contrast to the expected total reward criterion, the expected discounted reward criterion accounts for the “time value of money” – that is, receiving a reward at some future epoch is worth less than receiving an identical reward now. This is the most widely used criterion in infinite horizon models.

Finite horizon models

Define the *expected discounted reward* of the reward sequence $\mathcal{R}_\infty = (R_1, R_2, \dots)$ by

$$E \left[\sum_{n=1}^N \lambda^{n-1} R_n \right] \quad (2.24)$$

where the expectation is respect to the distribution of the sequence of rewards and the quantity $0 \leq \lambda < 1$ is the *discount factor*. In finite horizon models, discounting has no impact on theory or algorithms, but may affect a decision maker's preference for policies. For example, with a discount factor close to 0, a decision maker will prefer actions that lead to larger immediate rewards, in contrast to a situation with a discount factor close to 1 when future rewards would be of greater significance.

Definition 2.4. For all $s \in S$, the *finite horizon expected discounted reward* of policy π is

$$v_\lambda^\pi(s) := E_s^\pi \left[\sum_{n=1}^N \lambda^{n-1} R_n \right]. \quad (2.25)$$

This definition is analogous to the expected total reward case. Above, the expectation is respect to the probability distribution of random rewards that result under different realizations of the stochastic process determined by π as described in Section 2.2.4.

Infinite horizon models

Discounting plays a fundamental role in the application and analysis of infinite horizon models and is the most studied optimality criterion.

Definition 2.5. For all $s \in S$, the *infinite horizon expected discounted reward* of policy π is

$$v_\lambda^\pi(s) := \lim_{N \rightarrow \infty} E_s^\pi \left[\sum_{n=1}^N \lambda^{n-1} R_n \right]. \quad (2.26)$$

Unlike in the expected total reward case, no special model structure is required for the infinite sum (2.26) to converge, only that $E[\lambda^{n-1} R_n]$ decays sufficiently quickly. To ensure this, it is assumed throughout the book that rewards are bounded. This means that there exists a finite M for which,

Bounded reward assumption:

$$|r(s, a, j)| \leq M \text{ for all } a \in A_s, s \in S \text{ and } j \in S. \quad (2.27)$$

Since the sum of a geometric series satisfies $\sum_{n=1}^{\infty} \lambda^{n-1} = 1/(1 - \lambda)$, under the bounded reward assumption it follows that

$$\frac{-M}{1 - \lambda} \leq v_{\lambda}^{\pi}(s) \leq \frac{M}{1 - \lambda} \quad (2.28)$$

for all $\pi \in \Pi^{\text{HR}}$. Equation (2.28) shows the key role of the assumption that $\lambda < 1$. As noted previously when $\lambda = 1$ the series may not converge even when rewards are bounded.

Note that in contrast to the infinite horizon expected total reward, the quantity

$$\lim_{N \rightarrow \infty} \sum_{n=1}^N \lambda^{n-1} R_n$$

exists in a discounted model²¹.

Hence it follows²² that

$$\lim_{N \rightarrow \infty} E_s^{\pi} \left[\sum_{n=1}^N \lambda^{n-1} R_n \right] = E_s^{\pi} \left[\lim_{N \rightarrow \infty} \sum_{n=1}^N \lambda^{n-1} R_n \right] = E_s^{\pi} \left[\sum_{n=1}^{\infty} \lambda^{n-1} R_n \right]. \quad (2.29)$$

where the expression on the right is a shorthand for the middle expression.

²¹This is true because when rewards are bounded, the “tail” of the sum can be made arbitrarily small for every trajectory.

²²The result follows from the bounded convergence theorem.

Example 2.1 (ctd.) Returning to Example 2.1, this example shows that when $\lambda < 1$, the limit of the sequence of partial sums for each of the discounted sequences of rewards exists and is finite.

For ξ_1 , the sequence of its discounted partial sums σ'_1 is

$$\sigma'_1 = (0, -\lambda, -\lambda + \lambda^2, -\lambda + \lambda^2 - \lambda^3, \dots)$$

and for ξ_2 , the sequence of partial sums is

$$\sigma'_2 = (0, \lambda, \lambda - \lambda^2, \lambda - \lambda^2 + \lambda^3, \dots).$$

Thus the limit for ξ_1 is given by^a

$$\sum_{n=1}^{\infty} \lambda^{n-1} R_n = -\frac{\lambda}{1+\lambda}.$$

Similarly, the limit of the sum of discounted rewards of sequence ξ_2 equals $\lambda/(1+\lambda)$. Thus when $\lambda < 1$

$$E \left[\lim_{N \rightarrow \infty} \sum_{n=1}^N \lambda^{n-1} R_n \mid X_1 = s_1 \right] = -0.5 \frac{\lambda}{1+\lambda} + 0.5 \frac{\lambda}{1+\lambda} = 0.$$

Since $\sigma'_1 + \sigma'_2 = (0, 0, 0, \dots)$, it follows that

$$\lim_{N \rightarrow \infty} E \left[\sum_{n=1}^N \lambda^{n-1} R_n \mid X_1 = s_1 \right] = 0,$$

since the expectation equals 0 for any finite N . Hence (2.29) holds in this model.

^aTo see this, let $X = 1 - \lambda + \lambda^2 + \dots$. Then $X = 1 - \lambda X$ so that $X = 1/(1+\lambda)$. Since the limit for $\xi_1 = X - 1$, the result follows.

An alternative and important interpretation of discounting

Discounting arises naturally in models where rewards are units of currency. Due to inflation, near-term rewards are preferable to future rewards. Discounting also arises in another, rather surprising way which extends its applicability.

Suppose the planning horizon length is not fixed but represented by a random variable, T , distributed according to a geometric distribution²³ with parameter λ , independent of (R_1, R_2, \dots) or $(X_1, Y_1, X_2, Y_2, \dots)$.

²³A random variable T follows a *geometric distribution* with parameter λ if $P[T = k] = (1-\lambda)\lambda^{k-1}$ for $k = 1, 2, \dots$

Let the expected total reward over a random horizon of length T be

$$E_T \left[E \left[\sum_{n=1}^T R_n \right] \right], \quad (2.30)$$

where E_T denotes the expectation with respect to T . Assuming R_n is bounded and $0 \leq \lambda < 1$,

$$\begin{aligned} E_T \left[E \left[\sum_{n=1}^T R_n \right] \right] &= E \left[\sum_{k=1}^{\infty} (1-\lambda) \lambda^{k-1} \sum_{n=1}^k R_n \right] \\ &= E \left[\sum_{n=1}^{\infty} \sum_{k=n}^{\infty} (1-\lambda) \lambda^{k-1} R_n \right] = E \left[\sum_{n=1}^{\infty} \lambda^{n-1} R_n \right] \end{aligned}$$

where the last inequality follows from the relationship $\sum_{k=n}^{\infty} \lambda^{k-1} = \lambda^{n-1}/(1-\lambda)$. The assumptions on R_n and λ allow interchange of the order of summation. Thus

$$E_T \left[E \left[\sum_{n=1}^T R_n \right] \right] = E \left[\sum_{n=1}^{\infty} \lambda^{n-1} R_n \right]. \quad (2.31)$$

When the rewards are generated by policy π starting from state s , this representation of the expected discounted reward follows:

Alternative representation for the expected discounted reward:

$$v_{\lambda}^{\pi}(s) = E_T \left[E_s^{\pi} \left[\sum_{n=1}^T R_n \right] \right]. \quad (2.32)$$

This result provides an alternative interpretation for discounting. Since the random variable T may be regarded as the time until the first “failure” in an independent, identically distributed series of Bernoulli trials with “success” probability λ , it means that when a decision maker uses expected total reward to evaluate policies in a system that terminates at the time of a random failure independent of the decision maker’s policy it is equivalent to using the expected discounted reward. This is particularly appropriate in applications where the process can terminate suddenly for exogenous reasons. Examples include:

- **Animal behavior modeling**, when the animal might die of unanticipated causes (e.g., predation) independent of its decision making. See Section 3.5.
- **Clinical decision making**, in which a patient may die as a consequence of some event independent of the disease being treated. See Section 3.6.

Thus, discounting makes sense even in models that do not use economic values for rewards.

Note this provides an approach for estimating a discounted reward through simulation. The two steps are:

1. Generate T from a geometric distribution.
2. Apply Algorithm 2.1

Optimal policies

Definition 2.6. A policy π^* is *discount optimal* if

$$v_\lambda^{\pi^*}(s) \geq v_\lambda^\pi(s) \quad (2.33)$$

for all $s \in S$ and $\pi \in \Pi^{\text{HR}}$.

Chapter 5 analyzes infinite horizon Markov decision process models under this optimality criterion.

2.3.3 Long-run average reward for an infinite horizon model

In contrast to discounting, which emphasizes short term behavior, the long-run average reward criterion focuses on steady state or limiting behavior of derived stochastic processes. For that reason, the long-run average reward is most appropriate for infinite horizon, non-terminating models with frequent decision epochs. Note that in finite horizon models, the average reward is equivalent to the expected total reward (why?).

As an example, consider a queuing system in which the decision maker inspects the system state very frequently, such as every second, and decides which service rate to use. Section 3.4 provides a rigorous formulation of such a problem. For such a problem:

- The **expected total reward** criterion would not be able to distinguish between policies because expected total costs or rewards would be unbounded.
- The **expected discounted reward** criterion would not be appropriate because the decision maker is interested in long-term system performance. Given the time scale of decision making, rewards at future decision epochs should **not** be less valuable than current rewards. However, if random failure of the system could occur discounting may be appropriate but it would require discount rates very close to 1.

The *average reward*²⁴ or more accurately the *long-run average reward* of the reward

²⁴This quantity is interchangeably referred to as the *long-run average reward* or *gain*.

sequence (R_1, R_2, \dots) is given by

$$\lim_{N \rightarrow \infty} \frac{1}{N} E \left[\sum_{n=1}^N R_n \right].$$

When rewards are generated by a policy π , the average reward is defined as follows.

Definition 2.7. For all $s \in S$, the *average reward*, $g^\pi(s)$, of policy π is

$$g^\pi(s) := \lim_{N \rightarrow \infty} \frac{1}{N} E_s^\pi \left[\sum_{n=1}^N R_n \right]. \quad (2.34)$$

Note that the limit in (2.34) exists for all stationary policies when S is finite, as will be shown in Chapter 7. It need not exist when S is countable or policies are history-dependent. In such cases replace the limit above by the “lim inf” or the “lim sup” (see Chapter 7). The quantity g^π is sometimes referred to as the *gain* because the expected total reward in state s grows (“gains value”) at rate $g^\pi(s)$ per epoch in the limit.

Similarly to the discounted model, when rewards are bounded,

$$\lim_{N \rightarrow \infty} \frac{1}{N} E_s^\pi \left[\sum_{n=1}^N R_n \right] = E_s^\pi \left[\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N R_n \right]. \quad (2.35)$$

Example 2.1 (ctd.) Returning to Example 2.1, this example shows that when rewards are bounded, the limit of the sequence of the average partial sums for each of the sequences of rewards exists and is finite.

For ξ_1 , the sequence of average partial sums σ_1'' is

$$\sigma_1'' = \left(0, -\frac{1}{2}, 0, -\frac{1}{4}, \dots \right)$$

so that its limit exists and equals 0. Similarly, the limit of average of partial sums for ξ_2 equals 0. Thus

$$E \left[\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N R_n \middle| X_1 = s_1 \right] = 0.$$

since the expectation equals 0 for any finite N . Hence it is easy to see that (2.35) holds in this model.

Optimal policies

Definition 2.8. A policy π^* is *average reward optimal* or *gain optimal* if

$$g^{\pi^*}(s) \geq g^\pi(s) \quad (2.36)$$

for all $\pi \in \Pi^{\text{HR}}$ and $s \in S$.

Chapter 7 analyzes infinite horizon Markov decision process models under this optimality criterion.

2.4 The one-period problem: A fundamental building block

One-period models are the building blocks of Markov decision processes. Most of the algorithms presented later in the book are based on decomposing a multi-period model into a series of interlinked one-period models.

A one-period model begins with a decision epoch, evolves for one period of time and ends at the terminal non-decision epoch. Thus, it is the simplest representation of a finite horizon Markov decision process, namely the case of $N = 2$. In it, a policy is a single Markovian decision rule and the derived stochastic process is (X_1, Y_1, X_2) . Note that in a one-period model, history-dependent policies and Markovian policies are equivalent because the only information available to the decision maker at the first (only) decision epoch is the state.

Figure 2.6 illustrates the timing of events and the nomenclature of the one-period problem.

Let S denote the state space and A_s denote the sets of actions defined for each $s \in S$. Assume that S and each A_s are discrete and finite. Let $r_1(s, a, j)$ denote the reward function in period 1, $p_1(j|s, a)$ denote the transition probability function in period 1, and $r_2(j)$ denote the terminal reward function. Analysis of this simple model provides intuition for more complex models.

Definition 2.9. Given a policy π in a one-period problem, the *policy value function* denoted $v^\pi(s)$, is the expected total reward when the process starts in state s at decision epoch 1. It is given by

$$v^\pi(s) := E^\pi[r_1(X_1, Y_1, X_2) + r_2(X_2)|X_1 = s], \quad (2.37)$$

where the expectation is with respect to the probability distribution of (X_1, Y_1, X_2) induced by policy π when the system starts in state $s \in S$.

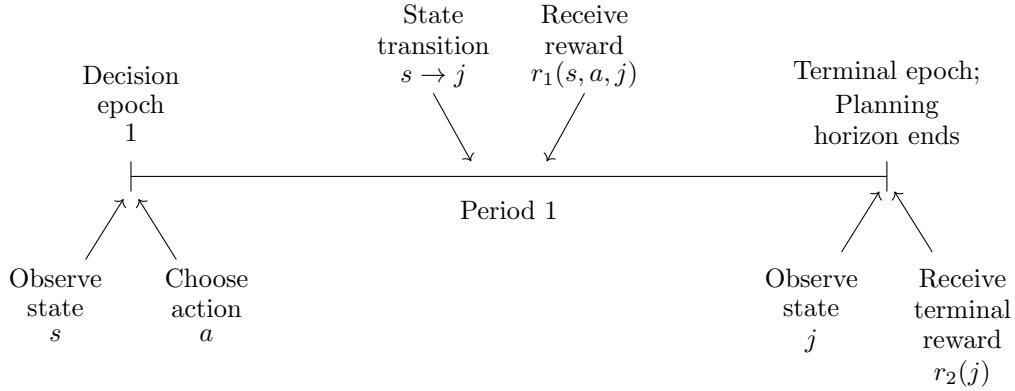


Figure 2.6: Illustration of timing of events and notation for the one-period problem. Action a may be chosen deterministically or probabilistically. Technically speaking, there is no need to observe state j since it is only used to evaluate the terminal reward $r_2(j)$.

Suppose $\pi = (d_1)$ is deterministic and $d_1(s) = a'$. Since s is fixed and $d_1(s) = a'$, the only random quantity in (2.37) is X_2 . Under this assumption, (2.37) is equivalent to

$$v^\pi(s) = \sum_{j \in S} p_1(j|s, a')(r_1(s, a', j) + r_2(j)). \quad (2.38)$$

Therefore, computing $v^\pi(s)$ requires only $p_1(j|s, a') = P[X_2 = j|X_1 = s, Y_1 = a']$; the distribution of $r_1(X_1, Y_1, X_2) + r_2(X_2)$ does not need to be explicitly derived. Consequently, this avoids an enumeration of all of the latter's possible values as described in Section 2.2.4. Expressions of the form (2.38) appear frequently in the book, so it is important to understand why (2.38) is equivalent to (2.37).

When d_1 is randomized, Y_1 is also random with $P[Y_1 = a|X_1 = s] = w_{d_1}^1(a|s)$. In this case

$$v^\pi(s) = \sum_{a' \in A_s} w_{d_1}^1(a'|s) \sum_{j \in S} p_1(j|s, a')(r_1(s, a', j) + r_2(j)). \quad (2.39)$$

In both cases, $v^\pi(s)$ is referred to as the *value* of policy π .

2.4.1 Optimal value functions

The following definition is fundamental.

Definition 2.10. The *optimal value function* in a finite horizon Markov decision process, $v^*(s)$, satisfies

$$v^*(s) := \sup_{\pi \in \Pi^{\text{HR}}} v^\pi(s) \quad (2.40)$$

for all $s \in S$.

Sometimes, the optimal value function is simply referred to as the *value function*. The adjective “optimal” distinguishes it from a policy value function.

In the one-period model, $\Pi^{\text{HR}} = \Pi^{\text{MR}}$. This identity holds because in a one-period model, the only decision is at the first decision epoch and the history at that decision epoch is the state s . Thus in the one-period model,

$$v^*(s) = \sup_{\pi \in \Pi^{\text{MR}}} v^\pi(s), \quad (2.41)$$

where $v^\pi(s)$ is given by (2.39). As a result of (2.39) and (2.41)

$$v^*(s) = \sup_{w \in \mathcal{P}(A_s)} \left\{ \sum_{a \in A_s} w(a) \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\}, \quad (2.42)$$

where $\mathcal{P}(A_s)$ denotes the set of probability distributions on A_s . The reason that (2.42) holds is that there is a one-to-one relationship between the set $\mathcal{P}(A_s)$ and the set of Markovian randomized decision rules. Since $\mathcal{P}(A_s)$ is an uncountable infinite set²⁵, \sup ²⁶ instead of \max appears in (2.42).

The following easily proved lemma is fundamental here and in latter chapters.

Lemma 2.1. Let U be an arbitrary finite set, let $f(\cdot)$ be a real-valued function on U , $\mathcal{P}(U)$ denote the set of probability distributions on U and $w \in \mathcal{P}(U)$. Then

$$\max_{u \in U} f(u) \geq \sum_{u \in U} w(u) f(u) \quad (2.43)$$

and

$$\max_{u \in U} f(u) = \sup_{w \in \mathcal{P}(U)} \sum_{u \in U} w(u) f(u). \quad (2.44)$$

Proof. To prove (2.43)

$$\sum_{u \in U} w(u) f(u) \leq \sum_{u \in U} w(u) \left(\max_{u \in U} f(u) \right) = \max_{u \in U} f(u).$$

²⁵For each state, the set of randomized decision rules in that state can be represented by a unit simplex corresponding to *all* probability distributions on A_s . The set Π^{MR} is equivalent to the Cartesian product of these simplices.

²⁶When maximizing over a non-finite set, “sup” is used even when it is attained to emphasize that the set is not finite. This point is explored further in the book appendix.

To prove (2.44), choose $u^* \in \arg \max_{u \in U} f(u)$ and define $w^* \in \mathcal{P}(U)$ by $w^*(u) = 1$ if $u = u^*$ and $w^*(u) = 0$, if $u \neq u^*$. Then, by the definition of w^* and (2.43),

$$\sum_{u \in U} w^*(u) f(u) = \max_{u \in U} f(u) \geq \sum_{u \in U} w(u) f(u)$$

for all $w \in \mathcal{P}(U)$ so the “sup” is attained by w^* and the result follows. \square

Applying Lemma 2.1 to (2.42) gives the first key result for one-period models.

Theorem 2.1. The optimal value function in a one-period model satisfies

$$v^*(s) = \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\} \quad (2.45)$$

for all $s \in S$.

Proof. Choose $s \in S$, set $U = A_s$ and

$$f(s) = \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)),$$

and apply (2.44). \square

The main consequence of this result is that in a one-period problem the value can be obtained by restricting attention to Markovian deterministic policies, which are equivalent to Markovian deterministic decision rules. The following corollary states this result more formally.

Corollary 2.1. In a one-period model,

$$v^*(s) = \max_{\pi \in \Pi^{\text{MD}}} v^\pi(s) \quad (2.46)$$

for all $s \in S$.

The result is expressed this way so that it generalizes to multi-period models.

2.4.2 Optimal policies

Attention now turns to identifying optimal policies. The concept of an “arg max” is critical and will be used throughout the book.

Definition 2.11. Let U denote an arbitrary set and $f(u)$ denote a real-valued function on U that attains its maximum on U . The *arg max* of $f(u)$ is defined by

$$\arg \max_{u \in U} f(u) := \left\{ u^* \in U \mid f(u^*) = \max_{u \in U} f(u) \right\}. \quad (2.47)$$

The “arg max” of a function returns the set of arguments that maximize the function. Thus, in general, the arg max is a set with multiple elements. Only when there is a unique maximizer of a function does the arg max return a single value. Of course the definition assumes that the maximum is attained. If not, the arg max is empty.

The following theorem shows how to identify optimal policies in the one-state model.

Theorem 2.2. In a finite state and action one-period Markov decision problem, let $d_1^*(\cdot)$ denote a deterministic decision rule that for each $s \in S$ satisfies

$$d_1^*(s) \in \arg \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\}. \quad (2.48)$$

Then the policy $\pi^* = (d_1^*)$ is an optimal policy.

Proof. From (2.48) and (2.38),

$$v^{\pi^*}(s) = \sum_{j \in S} p_1(j|s, d_1^*(s))(r_1(s, d_1^*(s), j) + r_2(j)) \quad (2.49)$$

$$= \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\} = v^*(s) \quad (2.50)$$

for all $s \in S$. Hence from (2.41), for all $s \in S$, $v^{\pi^*}(s) \geq v^\pi(s)$ for all $\pi \in \Pi^{\text{MR}}$ so that it is an optimal policy. \square

Some observations follow:

1. **Greedy actions:** Any action that achieves the maximum in the right hand side of (2.48) is referred to as a *greedy* action.
2. **Independent problems:** To find an optimal policy in this model, one solves an independent problem (2.48) for each state $s \in S$.
3. **Inter-temporal trade-offs:** Expressions like

$$\max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\}$$

appear frequently in Markov decision process theory and algorithms. This maximization highlights the trade-off between the immediate reward $r_1(s, a, j)$ and the future reward $r_2(j)$. This notion of balancing a current or *myopic* reward with future rewards is fundamental to Markov decision processes.

4. **Future values:** The future reward is encapsulated in the model component $r_2(X_2)$ in the one-period model. An important question in multi-period finite or infinite horizon models is: *Can the future reward be replaced with a single function that captures the cumulative reward associated with system evolution beyond the current decision epoch?* The answer is “yes”. Subsequent chapters elaborate on this key concept.
5. **Rollout and approximations:** In modern applications with very large state spaces, the terminal reward might represent an approximation to the “value” of being in state s that has been determined “offline”. Then to determine the next action in state s “online”, the decision maker solves the one-period problem and uses the greedy action. Such an approach has been used to determine good strategies in checkers, chess, go and backgammon.

2.4.3 State-action value functions

A significant development in computer-based model-free reinforcement learning is to shift the focus from value functions to *state-action value functions*. While not essential here and in Chapters 4–8, they play a key role in Chapters 10 and 11.

In general, a state-action value function gives the expected total reward (or discounted reward) when choosing action a in state s . For problems with planning horizons exceeding $N = 2$, the state-action value function must be distinguished according to whether the decision maker follows a specific policy or the optimal policy after action specification. This distinction is irrelevant in one-period problems in which there are no decisions after the first decision epoch.

Definition 2.12. In a one-period problem, the *state-action value function*, $q(s, a)$, is defined for all $a \in A_s$ and $s \in S$ by

$$q(s, a) := \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)). \quad (2.51)$$

State-action value functions often referred to as q -functions.

The state-action value function provides the decision maker with all the information necessary to find policy value functions, optimal policies and optimal value functions. Let $\pi = (d_1)$ denote an arbitrary deterministic policy. Then from (2.38), its value is given by

$$v^\pi(s) = q(s, d_1(s)). \quad (2.52)$$

Moreover, as a direct result of Theorem 2.1,

$$v^*(s) = \max_{a \in A_s} q(s, a) \quad (2.53)$$

and from Theorem 2.2

$$d_1^*(s) \in \arg \max_{a \in A_s} q(s, a). \quad (2.54)$$

Finding $d_1^*(s)$ using (2.54) is equivalent to greedy action choice with the q -function.

Figure 2.6 provides another way to distinguish state-action value functions from optimal value functions. The optimal value function corresponds to the largest total reward that can be obtained when the system starts in state s *before* choosing an action at the first decision epoch, while the (optimal) state-action value function refers to the optimal total reward *after* choosing action a in state s . Since there are no further decisions after action choice at decision epoch 1, the state-action value function does not involve a maximization.

As covered in Part III of this book, one might choose to estimate $q(s, a)$ by simulation. In such a situation, (2.52) can be used to find the value of a policy and (2.54) and (2.53) to find an optimal policy and its value, respectively.

2.4.4 Summary of results for a one-period problem

Putting this all together, the following has been demonstrated for a one-period model:

1. There exists a Markovian deterministic policy $\pi^* = (d_1^*)$ that is optimal in the class of all policies.
2. To “solve” the one-period problem under the expected total reward criterion, compute $v^*(s)$ for all $s \in S$ according to

$$v^*(s) = \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\} \quad (2.55)$$

and choose $d_1^*(s)$ so that

$$d_1^*(s) \in \arg \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\}. \quad (2.56)$$

3. The expected total reward of the optimal policy $\pi^* = (d_1^*)$, satisfies

$$\begin{aligned} v^{\pi^*}(s) &= \sum_{j \in S} p_1(j|s, d_1^*(s))(r_1(s, d_1^*(s), j) + r_2(j)) \\ &= \max_{a \in A_s} \left\{ \sum_{j \in S} p_1(j|s, a)(r_1(s, a, j) + r_2(j)) \right\} = v^*(s). \end{aligned}$$

The goal in later parts of the book will be to generalize these ideas to multi-period models under different optimality criteria.

2.5 A two-state model

This section introduces a simple example, frequently referred to as *the two-state model*, that illustrates model components and notation, and previews calculations that will be seen later in the book. It formulates a finite horizon version; an infinite horizon version is analyzed in depth in later chapters. To simplify notation assume the transition probabilities and rewards are stationary, that is, they do not vary from decision epoch to decision epoch. Moreover, the reward is a function of the current state, the current chosen action, and the subsequent state.

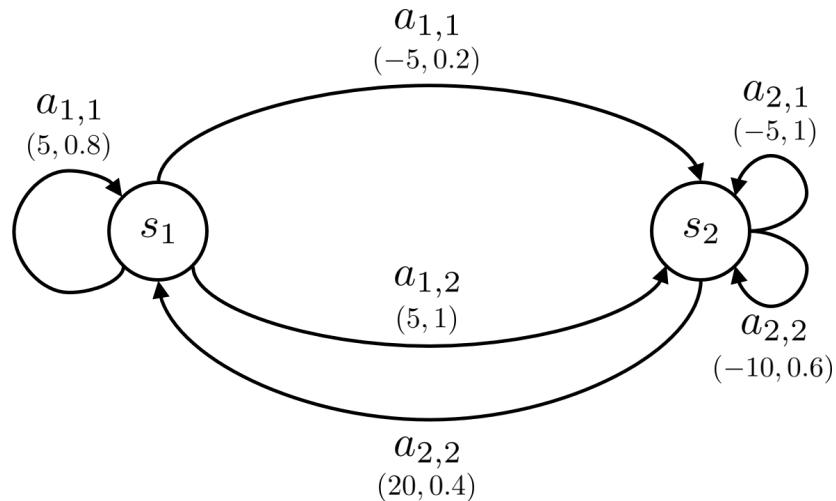


Figure 2.7: Graphical representation of two-state model. Circles denote states, arcs denote actions and transitions between states, and the expressions in parentheses denote rewards and transition probabilities, respectively. Zero probability transitions have been omitted.

Example 2.2. This example provides a concrete realization of the two-state model graphically represented in Figure 2.7. In this model, there are two states, and two actions to choose from in each state. Actions $a_{1,2}$ and $a_{2,1}$ result in deterministic outcomes – when the decision maker chooses these actions, transitions to a specified state occur with certainty. Consequently, rewards $r_n(s_1, a_{1,2}, s_1)$ and $r_n(s_2, a_{2,1}, s_1)$ are superfluous since they correspond to outcomes that cannot occur. Assume terminal rewards of 0. A precise formulation follows.

Decision epochs:

$$\{1, 2, \dots, N\}, \quad N < \infty$$

States:

$$S = \{s_1, s_2\}$$

Actions:

$$A_{s_1} = \{a_{1,1}, a_{1,2}\}, \quad A_{s_2} = \{a_{2,1}, a_{2,2}\}$$

Rewards: For $n = 1, 2, \dots, N$;

$$r_n(s_1, a_{1,1}, s_1) = 5, \quad r_n(s_1, a_{1,1}, s_2) = -5$$

$$r_n(s_1, a_{1,2}, s_1) = 0, \quad r_n(s_1, a_{1,2}, s_2) = 5$$

$$r_n(s_2, a_{2,1}, s_1) = 0, \quad r_n(s_2, a_{2,1}, s_2) = -5$$

$$r_n(s_2, a_{2,2}, s_1) = 20, \quad r_n(s_2, a_{2,2}, s_2) = -10$$

$$r_N(s_1) = 0, \quad r_N(s_2) = 0$$

Transition probabilities: For $n = 1, 2, \dots, N$;

$$p_n(s_1|s_1, a_{1,1}) = 0.8, \quad p_n(s_2|s_1, a_{1,1}) = 0.2$$

$$p_n(s_1|s_1, a_{1,2}) = 0, \quad p_n(s_2|s_1, a_{1,2}) = 1$$

$$p_n(s_1|s_2, a_{2,1}) = 0, \quad p_n(s_2|s_2, a_{2,1}) = 1$$

$$p_n(s_1|s_2, a_{2,2}) = 0.4, \quad p_n(s_2|s_2, a_{2,2}) = 0.6$$

2.5.1 A one-period version

This section shows how to compute optimal value functions and policies directly, and with state-action value functions as intermediaries, in Example 2.2. Recall that in a one-period model, $N = 2$.

Optimal value functions

From (2.45), $v^*(s)$ satisfies

$$\begin{aligned} v^*(s_1) &= \max_{a \in \{a_{1,1}, a_{1,2}\}} \{p_1(s_1|s_1, a)(r_1(s_1, a, s_1) + r_2(s_1)) + p_1(s_2|s_1, a)(r_1(s_1, a, s_2) + r_2(s_2))\} \\ &= \max\{0.8 \times (5 + 0) + 0.2 \times (-5 + 0), 5 + 0\} = \max\{3, 5\} = 5. \end{aligned}$$

and

$$v^*(s_2) = \max_{a \in \{a_{1,1}, a_{1,2}\}} \{p_1(s_1|s_1, a)(r_1(s_1, a, s_1) + r_2(s_1)) + p_1(s_2|s_1, a)(r_1(s_1, a, s_2) + r_2(s_2))\}$$

$$= \max\{-5 + 0, 0.4 \times (20 + 0) + 0.6 \times (-10 + 0)\} = \max\{-5, 2\} = 2.$$

The first term in the “max” corresponds to action $a_{i,1}$, $i = 1, 2$ and the second term in the “max” corresponds to $a_{i,2}$, $i = 1, 2$. Thus, from (2.48), $\pi^* = (d_1^*)$ where $d_1^*(s_1) = a_{1,2}$ and $d_1^*(s_2) = a_{2,2}$ and from (2.38), $v^{\pi^*}(s_1) = 5$ and $v^{\pi^*}(s_2) = -2$.

Observe that this calculation was particularly simple because the terminal reward in each state equals 0. Exercise 5 asks you to explore the sensitivity of the optimal decision to the terminal reward.

State-action value functions

These calculations are now repeated using state-action value functions. From (2.51),

$$\begin{aligned} q(s_1, a_{1,1}) &= p(s_1|s_1, a_{1,1})r(s_1, a_{1,1}, s_1) + p(s_2|s_1, a_{1,1})r(s_1, a_{1,1}, s_2) = 3 \\ q(s_1, a_{1,2}) &= p(s_1|s_1, a_{1,2})r(s_1, a_{1,2}, s_1) + p(s_2|s_1, a_{1,2})r(s_1, a_{1,2}, s_2) = 5 \\ q(s_2, a_{2,1}) &= p(s_1|s_2, a_{2,1})r(s_1, a_{2,1}, s_1) + p(s_2|s_2, a_{2,1})r(s_1, a_{2,1}, s_2) = -5 \\ q(s_2, a_{2,2}) &= p(s_1|s_2, a_{2,2})r(s_2, a_{2,2}, s_1) + p(s_2|s_2, a_{2,2})r(s_2, a_{2,2}, s_2) = 2. \end{aligned}$$

Hence from (2.54)

$$\begin{aligned} d_1^*(s_1) &= \arg \max_{a \in \{a_{1,1}, a_{1,2}\}} q(s_1, a) = a_{1,2} \\ d_1^*(s_2) &= \arg \max_{a \in \{a_{1,1}, a_{1,2}\}} q(s_2, a) = a_{2,2} \end{aligned}$$

and from (2.53)

$$\begin{aligned} v^*(s_1) &= \max_{a \in \{a_{1,1}, a_{1,2}\}} q(s_1, a) = 5 \\ v^*(s_2) &= \max_{a \in \{a_{1,1}, a_{1,2}\}} q(s_2, a) = 2 \end{aligned}$$

in agreement with direct calculation of optimal value functions.

As an aside, when coding algorithms in more complex problems, it was easier to organize calculations by first specifying $q(s, a)$.

2.5.2 Policies in a two-period version of the two-state model

In a one-period model, history-dependent and Markovian decision rules (and policies) are equivalent since the history at decision epoch 1 is the same as the state at decision epoch 1. Thus, to provide examples of the different types of policies, this section considers a two-period version. In a two-period example $N = 3$ representing two-decision epochs and a subsequent point in time when the terminal state is realized.

A Markovian deterministic policy in Π^{MD}

Let $\pi_1 = (d_1, d_2)$, where

$$d_1(s) = \begin{cases} a_{1,1}, & s = s_1 \\ a_{2,1}, & s = s_2 \end{cases} \quad \text{and} \quad d_2(s) = \begin{cases} a_{1,2}, & s = s_1 \\ a_{2,1}, & s = s_2. \end{cases} \quad (2.57)$$

In state s_1 , π_1 chooses action $a_{1,1}$ at the first decision epoch and $a_{1,2}$ at the second decision epoch. In state s_2 , π chooses action $a_{2,1}$ at both decision epochs.

The expected total reward generated by policy π_1 for each starting state was computed by enumerating sample paths and assigning probabilities to each. These simple calculations are summarized in Table 2.1.

In state s_1 π_1 chooses $a_{1,1}$ and remains in state s_1 with probability 0.8 and jumps to state s_2 with probability 0.2. The corresponding rewards for these transitions are 5 and -5 , respectively. At decision epoch 2, if the system is in state s_1 , policy π_1 chooses $a_{1,2}$ resulting in a transition to state s_2 with certainty and a reward of 5. This sample path, which occurs with probability 0.8, has a total reward of $5 + 5 = 10$.

If the state at decision epoch 2 is s_2 , this policy chooses action $a_{2,1}$, resulting in a self-transition and a reward of -5 . This sample path occurs with probability 0.2 and has a total reward of $-5 - 5 = -10$. Since the terminal rewards are 0, the expected total reward of this Markovian deterministic policy is $0.8(10) + 0.2(-10) = 6$. Note that in these calculations the state at the end of period 2 was required to evaluate the reward in period 2.

Starting state	Sample path	Total reward	Probability	Expected total reward
s_1	$(s_1, a_{1,1}, s_1, a_{1,2}, s_2)$	10	0.8	6
	$(s_1, a_{1,1}, s_2, a_{2,1}, s_2)$	-10	0.2	
s_2	$(s_1, a_{2,1}, s_2, a_{2,1}, s_2)$	-10	1	-10

Table 2.1: Evaluation of expected total reward for Markovian deterministic policy π_1 by enumerating sample paths.

Clearly these calculations are tedious when N and the states and action sets were larger. Inductive methods in Chapter 4 provide a more efficient way to evaluate the expected rewards in a finite horizon model.

If instead, the initial state is chosen *randomly* with $\rho(s_1) = p$ and $\rho(s_2) = 1 - p$ for some $0 \leq p \leq 1$, the expected total reward corresponding to π_1 equals $16p - 10$.

A Markovian randomized policy in Π^{MR}

Next is an example of a Markovian randomized policy $\pi_2 = (d_1, d_2)$. The randomized decision rule d_n , for $n = 1, 2$, choose action $a_{1,1}$ in state s_1 , with probability q_1^n and action $a_{1,2}$ with probability $1 - q_1^n$ and in state s_2 , choose action $a_{2,1}$ with probability

q_2^n and action $a_{2,2}$ with probability $1 - q_2^n$. Then, in the notation of Section 2.2.1, for $n = 1, 2$;

$$w_{d_n}^n(a|s) = \begin{cases} q_1^n, & a = a_{1,1}, s = s_1 \\ 1 - q_1^n, & a = a_{1,2}, s = s_1 \\ q_2^n, & a = a_{2,1}, s = s_2 \\ 1 - q_2^n, & a = a_{2,2}, s = s_2. \end{cases} \quad (2.58)$$

Computing the expected total reward of a Markovian randomized policy is slightly more complicated than in the Markovian deterministic case because it requires taking into account action choice probabilities specified by d_1 and d_2 . Table 2.2 exhibits these calculations when the process starts in state s_1 and $q_1^1 = 0.1$, $q_1^2 = 0.5$ and $q_2^2 = 0.7$. Note that since the system starts in s_1 with certainty, q_2^1 is not used in the calculations below. In this table, a sample path represents a realization of the stochastic process $(X_1, Y_1, X_2, Y_2, X_3)$. As above, the realization of X_3 is required to compute the reward in period 2 which depends on the state at end of the horizon s^3 .

Sample path	Total reward	Probability
$(s_1, a_{1,1}, s_1, a_{1,1}, s_1)$	$5 + 5 = 10$	$0.1 \times 0.8 \times 0.5 \times 0.8 = 0.0320$
$(s_1, a_{1,1}, s_1, a_{1,1}, s_2)$	$5 - 5 = 0$	$0.1 \times 0.8 \times 0.5 \times 0.2 = 0.0080$
$(s_1, a_{1,1}, s_1, a_{1,2}, s_2)$	$5 + 5 = 10$	$0.1 \times 0.8 \times 0.5 \times 1.0 = 0.0400$
$(s_1, a_{1,1}, s_2, a_{2,1}, s_2)$	$-5 - 5 = -10$	$0.1 \times 0.2 \times 0.7 \times 1.0 = 0.0140$
$(s_1, a_{1,1}, s_2, a_{2,2}, s_1)$	$-5 + 20 = 15$	$0.1 \times 0.2 \times 0.3 \times 0.4 = 0.0024$
$(s_1, a_{1,1}, s_2, a_{2,2}, s_2)$	$-5 - 10 = -15$	$0.1 \times 0.2 \times 0.3 \times 0.6 = 0.0036$
$(s_1, a_{1,2}, s_2, a_{2,1}, s_2)$	$5 - 5 = 0$	$0.9 \times 1.0 \times 0.7 \times 1.0 = 0.6300$
$(s_1, a_{1,2}, s_2, a_{2,2}, s_1)$	$5 + 20 = 25$	$0.9 \times 1.0 \times 0.3 \times 0.4 = 0.1080$
$(s_1, a_{1,2}, s_2, a_{2,2}, s_2)$	$5 - 10 = -5$	$0.9 \times 1.0 \times 0.3 \times 0.6 = 0.1620$

Table 2.2: Sample paths, rewards and probabilities corresponding to Markovian randomized policy π_2 starting in state s_1 . Zero-probability sample paths are not shown.

The following expression mathematically describes the quantities used to obtain the probability in the first row of Table 2.2.

$$\begin{aligned} P^{\pi_2}[(X_1, Y_1, X_2, Y_2, X_3) = (s_1, a_{1,1}, s_1, a_{1,1}, s_1)] \\ = w_{d_1}^1(a_{1,1}|s_1)p(s_1|a_{1,1}, s_1)w_{d_2}^2(a_{1,1}|s_1)p(s_1|a_{1,1}, s_1) \end{aligned}$$

Combining sample paths which generate the same reward provides the following probability distribution of total rewards generated by π_2 and enables computation of its mean and standard deviation.

From Table 2.3 it follows that the expected total reward corresponding to starting π_2 in s_1 is 2.45 and its standard deviation is 9.35.

Reward	Probability
25	0.1080
15	0.0024
10	0.0720
0	0.6380
-5	0.1620
-10	0.0140
-15	0.0036

Table 2.3: Derived probability distribution of total reward over decision epochs 1 and 2 for Markovian randomized policy π_2 starting in state s_1 .

A history-dependent deterministic policy in Π^{HD}

Let $\pi_3 = (d_1, d_2)$, where

$$d_1(h) = \begin{cases} a_{1,1}, & h = s_1 \\ a_{2,1}, & h = s_2 \end{cases} \quad \text{and} \quad d_2(h) = \begin{cases} a_{1,2}, & h = (s_1, a_{1,1}, s_1) \\ a_{2,1}, & h = (s_1, a_{1,1}, s_2) \\ a_{2,2}, & h = (s_2, a_{2,1}, s_2) \end{cases} \quad (2.59)$$

For this policy, the chosen action (at decision epoch 2) depends on the entire history. At decision epoch 1, the history is simply the initial state. But at decision epoch epoch 2, the history includes the initial state, the first action chosen, and the subsequent state. If the state is s_1 at decision epoch 2, there is only one way for this to happen, given d_1 (initial state s_1 with action $a_{1,1}$). However, there are two histories that lead to s_2 at epoch 2. Thus, the action chosen in state s_2 at epoch 2 depends on whether the process started in s_1 or s_2 : starting in s_1 leads to choosing action $a_{2,1}$ in the second decision epoch, otherwise $a_{2,2}$. Notice that if the policy selects action $a_{2,1}$ for the history $(s_2, a_{2,1}, s_2)$, then this policy coincides with the Markovian deterministic policy described previously.

It is left as an exercise to derive the probabilities of each realization and expected total reward.

A history-dependent randomized policy in Π^{HR}

Define the policy $\pi_4 = (d_1, d_2)$ as follows. At the first decision epoch, the randomized decision rule is described by the following probability distribution:

$$w_{d_1}^1(a|h) = \begin{cases} q_1^1, & a = a_{1,1}, h = s_1 \\ 1 - q_1^1, & a = a_{1,2}, h = s_1 \\ q_2^1, & a = a_{2,1}, h = s_2 \\ 1 - q_2^1, & a = a_{2,2}, h = s_2. \end{cases} \quad (2.60)$$

In state s_1 , action $a_{1,1}$ is chosen with probability q_1^1 and action $a_{1,2}$ is chosen with probability $1 - q_1^1$. In state s_2 , action $a_{2,1}$ is chosen with probability q_2^1 and action $a_{2,2}$ is chosen with probability $1 - q_{2,1}$.

At the second decision epoch, there are eight histories to consider, since there are two states and two actions in each state. However, $a_{1,2}$ and $a_{2,1}$ result in deterministic transitions to state s_2 , so two of the eight histories are impossible. The remaining six are listed below:

$$\begin{aligned} h_1 &:= (s_1, a_{1,1}, s_1) \\ h_2 &:= (s_1, a_{1,1}, s_2) \\ h_3 &:= (s_1, a_{1,2}, s_2) \\ h_4 &:= (s_2, a_{2,1}, s_2) \\ h_5 &:= (s_2, a_{2,2}, s_1) \\ h_6 &:= (s_2, a_{2,2}, s_2). \end{aligned}$$

Given these possible histories at the second decision epoch, a possible specification for a randomized decision rule is

$$w_{d_2}^2(a|h) = \begin{cases} q_i^2, & a = a_{1,1}, h = h_i \\ 1 - q_i^2, & a = a_{1,2}, h = h_i \end{cases} \quad (2.61)$$

for $i \in \{1, 5\}$ and

$$w_{d_2}^2(a|h) = \begin{cases} q_i^2, & a = a_{2,1}, h = h_i \\ 1 - q_i^2, & a = a_{2,2}, h = h_i \end{cases} \quad (2.62)$$

for $i \in \{2, 3, 4, 6\}$.

Again it is left as an exercise to derive the probabilities of each realization and expected total reward.

A stationary deterministic policy $\pi \in \Pi^{\text{SD}}$

Suppose $\pi_5 = (d, d)$, where

$$d(s) = \begin{cases} a_{1,1}, & s = s_1 \\ a_{2,1}, & s = s_2 \end{cases} \quad (2.63)$$

This policy uses the same decision rule in both decision epochs. It is similar to the Markovian deterministic policy above, except that if the system is in state s_1 at epoch 2, the SD policy chooses action $a_{1,1}$, whereas the MD policy chooses $a_{1,2}$.

A stationary randomized policy $\pi \in \Pi^{\text{SR}}$

The Markovian randomized policy above can be transformed into a stationary randomized policy $\pi = (d, d)$ by simply omitting its dependence on n , for example:

$$w_d(a|s) = \begin{cases} q_1, & a = a_{1,1}, s = s_1 \\ 1 - q_1, & a = a_{1,2}, s = s_1 \\ q_2, & a = a_{2,1}, s = s_2 \\ 1 - q_2, & a = a_{2,2}, s = s_2. \end{cases} \quad (2.64)$$

This policy has a stationary probability distribution of choosing action $a_{1,1}$ versus $a_{1,2}$ when the system is in state s_1 , and similarly when the state is s_2 .

Bibliographic remarks

See the historical summary in Section 1.1.1.

Exercises

1. Consider the following three-state model with $S = \{s_1, s_2, s_3\}$, actions $A_{s_i} = \{a_{i,1}, a_{i,2}\}$ for $i = 1, 2, 3$, rewards $r_n(s_i, a_{i,1}, s_j) = i - j$, $r_n(s_i, a_{i,2}, s_j) = j - i$ for $n = 1, 2, \dots, N - 1$ and $r_N(s_i) = -i^3$. Transition probabilities are given by $p_n(s_i|s_i, a_{i,1}) = 1 - 1/i$, $p_n(s_j|s_i, a_{i,1}) = 1/(2i)$ for $j \neq i$ and $p_n(s_i|s_i, a_{i,2}) = 1 - 1/(i + 1)$, $p_n(s_j|s_i, a_{i,2}) = 1/(2(i + 1))$ for $j \neq i$ for $n = 1, 2, \dots, N$.
 - (a) Provide a graphical representation of the model as in Figure 2.7.
 - (b) Represent a one- and two-period version of the model as a decision tree when $\rho(s_i) = P[X_1 = s_i] = 1/3$ for $i = 1, 2, 3$. Compute the expected total reward by rolling back calculations on each path through the tree.
 - (c) In a one-period version, find the distribution of X_2 for each possible initial state s_1, s_2, s_3 when the deterministic decision rule $d(s_i) = a_{i,1}$ is used at decision epoch 1.
 - (d) Use (2.38) to find the expected total reward $v^\pi(s)$ for each $s \in S$ in a one-period version for the policy π that uses the above decision rule d at decision epoch 1.
 - (e) In a one-period version, find the distribution of X_2 for each possible initial state s_1, s_2, s_3 when the randomized decision rule d' with distribution $P[d'(s_i) = a_{i,1}] = 1 - P[d'(s_i) = a_{i,2}] = e^{-0.5i}$ for $i = 1, 2, 3$ is used at decision epoch 1.
 - (f) Use (2.39) to find the expected total reward $v^{\pi'}(s)$ for each $s \in S$ in a one-period version for the policy π' that uses the above decision rule d' at decision epoch 1.
 - (g) Find the optimal policy in a one-period version by computing optimal value functions and state-action value functions.

2. Using 5,000 replications of Algorithm 2.1, simulate the total reward for the policies π and π' from Exercise 1 when $N = 2$.
 - (a) Estimate the expected total reward of each policy and compare your results to those in Exercise 1.
 - (b) Provide histograms of your estimates and comment on their shape and any differences you observe between the histograms of π and π' .
 - (c) Compute the standard deviation and 95th percentile of the total returns for each policy. Interpret these quantities verbally and note why one might be interested in such quantities.
 - (d) Suppose one measures the value of a reward stream (R_1, R_2) by multiplicative utility $e^{\gamma R_1} e^{\gamma R_2}$. Compare policies on the basis of their expected utility. Repeat parts (a) to (c) above using this utility function.
3. Show that when $\pi \in \Pi^{\text{MR}}$, the sequence of state and action pairs is a discrete time Markov Chain. Devise an example that shows that when $\pi \in \Pi^{\text{HR}}$, the sequence may not be Markov Chain.
4. Construct an example where an epoch-dependent state space S_n is appropriate (i.e., where using $S = \cup_n S_n$ would unnecessarily enlarge the state space at each decision epoch).
5. Write out (2.55) and (2.56) for the one-period version of the two-state problem in which $r_2(s_1) = c_1$ and $r_2(s_2) = c_2$. Plot the optimal policy as a function of the terminal rewards c_1 and c_2 .
6. Evaluate the probability distributions of stochastic process generated by the deterministic history-dependent policy π_3 and the randomized history-dependent policy π_4 for the two-period model in Section 2.5.2. To organize calculations create analogs of Tables 2.2 and 2.3.
7. *Consider the following deterministic model. Let $S = \{s_1, s_2\}$, $A_{s_1} = \{a_{1,1}, a_{1,2}\}$, $A_{s_2} = \{a_{2,1}, a_{2,2}\}$, $r(s_1, a_{1,1}) = r(s_1, a_{1,2}) = 2$, $r(s_2, a_{2,1}) = r(s_2, a_{2,2}) = -2$, and $p(s_1|s_1, a_{1,1}) = p(s_2|s_1, a_{1,2}) = p(s_1|s_2, a_{2,1}) = p(s_2|s_2, a_{2,2}) = 1$.
 - (a) Provide a graphical representation of the model as in Figure 2.7.
 - (b) Show that for each stationary policy π and each state $s \in S$ that the following limit exists

$$\lim_{N \rightarrow \infty} \frac{1}{N} E^\pi \left[\sum_{n=1}^N r(X_n, Y_n) \mid X_1 = s \right]. \quad (2.65)$$
 - (c) Consider a history-dependent policy π that when the initial state is s_1 chooses action $a_{1,1}$ for one period, then chooses $a_{1,2}$ so that the system

proceeds to s_2 and then chooses action $a_{2,2}$ so the system remains in s_2 for three periods, at which point it chooses action $a_{2,1}$ so that the system returns to s_1 and then chooses action $a_{1,1}$ so that it stays in state s_1 for $3^2 = 9$ periods and then chooses actions so that it proceeds to s_2 and remains there for $3^3 = 27$ periods and so forth.

Show that for π the limit in (2.65) does not exist by showing that

$$\liminf_{N \rightarrow \infty} \frac{1}{N} E^\pi \left[\sum_{n=1}^N r(X_n, Y_n) \mid X_1 = s \right] \neq \limsup_{N \rightarrow \infty} \frac{1}{N} E^\pi \left[\sum_{n=1}^N r(X_n, Y_n) \mid X_1 = s \right].$$

Chapter 3

Examples and Applications

This material will be published by Cambridge University Press as “Markov Decision Processes and Reinforcement Learning” by Martin L. Puterman and Timothy C. Y. Chan. This pre-publication version is free to view and download for personal use only. Not for re-distribution, re-sale, or use in derivative works. ©Martin L. Puterman and Timothy C. Y. Chan, 2025.

*For the things we have to learn before we can do them, we learn by doing
them.¹*

Aristotle, Greek philosopher, 384-322 BCE.

Representing a sequential decision problem as a Markov decision process requires specifying five objects: decision epochs, states, actions, transition probabilities and rewards. This chapter shows how to do so by identifying these objects in applications spanning a broad range of disciplines. The objective of this chapter is to highlight both the wide applicability of the MDP framework and the diverse aspects of model formulation. Note that several of these examples will be revisited in subsequent chapters.

We encourage readers to identify decision epochs, states, actions, transition probabilities and rewards on their own before looking at our approach. The chapter concludes with guidance on how to formulate Markov decision process models in general. The problems at the end of the chapter provide additional opportunities to formulate Markov decision processes or revise the provided examples when assumptions differ.

3.1 Inventory management

Inventory models represent some of the earliest and most widely studied applications of Markov decision processes. They concern determining appropriate inventory levels for

¹Aristotle [1925].

a product in the face of uncertain customer demand. The following passage from a New York Times article² nicely summarizes the core challenge of inventory management.

“When you have an inventory-based business, most people think only about the first order,” Mr. Green said. With long lead times from the factory in China, he was almost immediately trying to figure out how big his second and third orders should be. Underestimating would hurt not just sales but the overall status of his Amazon listing; overestimating would drain him of cash upfront, and he would incur further charges from Amazon for storing excess inventory in its warehouses.

As the quote indicates, having too little inventory results in losing sales and reputation, while having too much inventory leads to excess storage and capital charges. These costs are key components of inventory models; balancing this trade-off is the primary challenge.

The models considered here assume that an *inventory manager* periodically (hourly, daily, weekly or monthly) observes the inventory level of a product and, if deemed opportunistic, places a replenishment order with a *supplier*. This replenishment may arrive immediately, at the end of the current period (prior to the next review period) or several periods in the future. The delay between placing and receiving an order is referred to as a *lead time*.

Inventory management applications have many features that impact an Markov decision process formulation.

1. **Ordering costs** consist of a fixed component and a variable component. The fixed component K represents the administrative cost of placing an order³ and the variable component $c(u)$ represents the cost of ordering u units. Thus when an order of u units is placed, the cost is $K + c(u)$ and when no order is placed, the cost is 0.
2. **Holding costs**, denoted by $h(u)$, represent the cost to the inventory manager of storing u units of product for one period. This cost only applies when $u > 0$. It is convenient to assume that $h(0) = 0$.
3. **Lost sales** occur if there is insufficient inventory to fulfill demand in the current period.
4. **Backlogged demand** means that unmet demand may be fulfilled at some period in the future when inventory is available. *A different formulation applies depending on whether unmet demand is lost or backlogged.*

²Herrman [2021].

³In some applications, $K = 0$, but usually $K > 0$. It is the presence of this cost that makes it advantageous for the inventory manager to not order every period.

5. **Penalty costs**, denoted by $p(u)$, represent the cost to the inventory manager of backlogging u units of demand for one period. This cost applies only when $u < 0$. It is convenient to assume $p(0) = 0$.
6. **Total customer demand** in each period is random. It may arrive in one batch at a specific point in time within a period or may arrive at random times within a period. The timing of demand will impact the formulation. The demand distribution may be known or unknown, and may be static or time-varying. Let the non-negative⁴ random variable Z_n denote the total demand in period n . Its distribution will be represented by $P[Z_n = z]$.
7. **Revenue** corresponds to the amount spent by customers purchasing units of a product. It is represented by an increasing function $R(\cdot)$ and its relationship to demand depends on some of the above features. If unfulfilled demand is backlogged, it is simplest to assume that payments are made “up front”, that is when an order is placed⁵ so that if the demand is z , the revenue is $R(z)$. In the lost sales case, sales are capped by inventory on hand so that the revenue is $R(z)$ if the demand is less than inventory level s and $R(s)$ if the demand exceeds the inventory s . In this latter case, low inventory levels may result in lost sales and reduced revenue.
8. **Product shelf-life** may be finite or “infinite⁶”. Products with finite shelf lives are said to be *perishable*. Perishable items may last for one-period (newspapers, fresh bread, defrosted vaccine or a seat at a sporting event) or multiple-periods (blood products, consumer electronics or fashion goods). Non-perishable goods include tools, books and basic clothing. Our discrete-state formulation implicitly assumes product is only available in whole units.
9. **Scrap value**, denoted by $H(s)$, equals the value of the ending inventory when there are s units on hand at the end of a finite planning horizon. If $s \geq 0$, $H(s)$ may include any potential holding cost before the inventory is liquidated. If $s < 0$, then $H(s)$ represents a penalty associated with not being able to fulfill the $|s|$ items backlogged (e.g., the loss associated with buying these items at a premium from a back-up supplier, or a loss in goodwill due to lost sales).

3.1.1 Inventory management with backlogged demand

Formulating sequential inventory management as a Markov decision process requires precisely specifying the timing of events and assumptions. Figure 3.1 depicts the

⁴Most inventory models assume non-negative demand.

⁵Alternatively, payments may be made when the order is fulfilled.

⁶No product lasts forever. Assuming infinite shelf life is a modeling convention to facilitate representing future inventory levels as functions of the present inventory, the quantity ordered and the demand only.

event sequence for the following model. At the start of period N , the decision maker observes the inventory level s and decides on how many units a of product to order from the supplier to replenish inventory. Demand arrives randomly throughout the period and is fulfilled at the end of the period, using inventory on hand together with units that arrive from the current period's order. If demand exceeds inventory, it is backlogged for fulfillment from future inventory replenishments. An Markov decision process formulation follows.

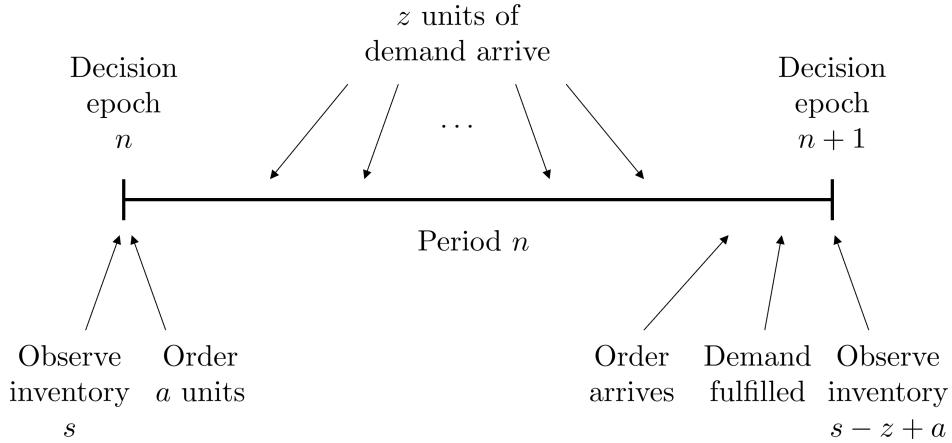


Figure 3.1: Timing of events in the periodic review inventory model.

Decision epochs: Decision epochs correspond to the times at which the decision maker reviews the inventory. This problem can be modeled either with a finite or infinite planning horizon. Hence

$$T = \{1, 2, \dots, N\}, \quad N \leq \infty.$$

States: States represent the number of units on hand at a decision epoch. A negative value indicates backlogged demand.

$$S = \{\dots, -2, -1, 0, 1, 2, \dots\}.$$

Note that the quantities backlogged and in inventory may be truncated at large values (e.g., the capacity of a warehouse) to ensure a finite state space. Doing so makes the formulation slightly more complex because of the resulting boundary conditions.

Actions: Actions represent the quantity ordered from the supplier for delivery prior to the next decision epoch. For each $s \in S$,

$$A_s = \{0, 1, 2, \dots\}.$$

Similar to the state space, the action space can be truncated so that it is finite.

Rewards: The reward in a period equals revenue minus cost. As noted above, revenue depends on whether unmet demand is backlogged or lost. Since this formulation assumes backlogging, when the total demand in a period is z , the revenue equals $R(z)$. Note that the inventory level at the next decision epoch, j is given by $j = s + a - z$ so that knowing only the inventory levels s and j and the order quantity a one can infer that the demand is given $z = s + a - j$.

There are three cost components: ordering costs, which are incurred only if $a > 0$, holding costs if s is positive and penalty costs if s is negative. A simplifying assumption is that holding and penalty costs are assessed at the beginning of the period based on the *starting* inventory.

The reward can be written⁷

$$r_n(s, a, j) = R(s + a - j) - K I_{\mathbb{Z}_{>0}}(a) - c(a) - h(s) I_{\mathbb{Z}_{>0}}(s) - p(-s) I_{\mathbb{Z}_{>0}}(s)$$

when $n < N$. The argument of $p(\cdot)$ is $-s$, which equals the backlogged demand when s is negative.

If N is finite, $r_N(s) = H(s)$, corresponding to the scrap value of the remaining inventory at the end of the planning horizon.

Transition probabilities: As noted above, when the state at decision epoch $n + 1$ is j , this means that demand in period n was $s + a - j$. Moreover, the assumption of non-negative demand ensures that the state at the next decision epoch cannot exceed $s + a$. Hence, the transition probabilities are

$$p_n(j|s, a) = \begin{cases} P[Z_n = s + a - j], & j \leq s + a, \text{ } j \text{ integer} \\ 0, & j = s + a + 1, s + a + 2, \dots \end{cases}$$

Application challenges

Applying this model presents several challenges. In particular, one must determine demand distributions, ordering costs, holding costs and penalty costs. Demand distributions may be estimated from historical data; parameterizing the model in terms of a known distribution reduces the challenge in estimating model parameters. Moreover it is likely that demand has seasonal components at the day, week and month levels.

Per unit ordering costs should be easily obtainable but fixed ordering costs may be more challenging to determine. The fixed component encompasses administrative, shipping and handling costs that may be hard to untangle from the per unit cost. Holding costs are real and involve cost of capital and space charges.

Penalty costs are the most challenging to determine since they involve intangibles such as loss of goodwill due to unfulfilled or delayed orders. Instead, the decision maker

⁷In the definition of the reward $\mathbb{Z}_{>0}$ denotes the positive integers so that $I_{\mathbb{Z}_{>0}}(x)$ equals 1 when x is positive and 0 otherwise.

may specify a service level such as 95% of orders be processed from stock on hand and use a constrained Markov decision process model formulation.

Note that major disruptions to supply chain or consumer behavior arising from, for example, a global pandemic, can lead to significant challenges in estimating appropriate parameters. Procurement costs might be much higher due to increased demand for raw materials and lower manufacturing capacity. Demand for certain products could be significantly increased or decreased compared to historical levels.

3.1.2 The newsvendor problem

This section describes a simple yet widely studied and applied inventory model referred to as the *newsvendor problem*⁸. This model is fundamental in the operations research literature and applies to perishable goods with a shelf-life of one period.

As shown in Figure 3.2, events unfold as follows. At the start of the period, the newsvendor purchases a units of product (newspapers) from a supplier at a cost of c per unit. After receiving the units, a random demand of z units arrives throughout the period. Items are sold at a price of g with $g > c$. Each unit sold results in a profit (selling price minus cost) $G := g - c$. At the end of the period each unsold item is disposed of at a scrap value of h where $c > h$ resulting in a loss (cost minus salvage value) $L := c - h$ and an ending of inventory zero.

The newsvendor seeks an order quantity a^* that maximizes the one-period expected reward. The choice of order quantity *trades off* ordering too many items and incurring a loss on unsold items with ordering too few items and foregoing potential additional profit. Another concern is that if a is set too low, the newsvendor will observe only *censored* demand complicating estimation of the demand distribution when it is unknown. Thus, initially the newsvendor may set a higher, trading off potential losses for gains in information about the demand distribution.

Formulation as a one-period Markov decision process follows.

Decision epochs: There is only one decision epoch so that $T = \{1\}$, consequently $N = 2$.

States: The formulation is a bit non-standard as the set of possible states differs at decision epoch 1 and at the terminal epoch. At the first decision epoch $S_1 = \{0\}$ corresponding to the initial inventory of 0. At the terminal epoch, $S_2 = A_0 = \{0, 1, \dots, a_{\max}\}$ where a_{\max} is defined below. This quantity is required to evaluate the final inventory and accounts for the implicit assumption that demand cannot be negative.

⁸This was first called a “newsboy” problem. It modeled the problem of a newsboy who purchased papers from a dealer, sold them at a street corner and returned the unsold papers to the dealer at the end of the day.

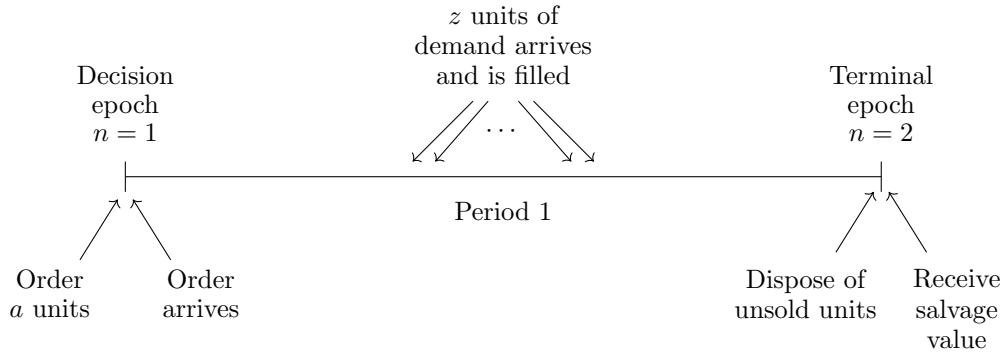


Figure 3.2: Timing of events in the newsvendor problem.

Actions: Actions represent the quantity ordered from the supplier for immediate delivery. This means that

$$A_0 = \{0, 1, 2, \dots, a_{\max}\},$$

where the set of actions is bounded by some “reasonable” value a_{\max} .

Rewards: The reward in the first period includes the cost of purchasing the product plus the revenue from sales prior to the end of the period. Thus

$$r_1(0, a, j) = \begin{cases} -ca + g(a - j) & \text{for } j \leq a \\ 0 & \text{for } j > a. \end{cases}$$

Note it is impossible for ending inventory to exceed a .

The terminal reward results from disposing of the unsold product. That is

$$r_2(s) = hs.$$

Note that on account of the transition probabilities, $s \leq a$.

Transition probabilities: The transition probabilities at decision epoch 1 account are computed differently depending on whether some or all inventory is sold. The case $j = 0$ occurs when the demand equals or exceeds the order quantity a . Moreover the inventory at the end of the period cannot exceed the order quantity.

$$p_1(j|s, a) = \begin{cases} P[Z = a - j] & \text{for } j = 1, \dots, a \\ \sum_{z=a}^{\infty} P[Z = z] & \text{for } j = 0 \\ 0 & \text{for } j > a. \end{cases}$$

Optimizing the order quantity

Suppose the newsvendor orders a units and the demand is z units. Then the total revenue combining $r_1(0, a, j)$ and $r_2(j)$ is given by

$$\begin{aligned} r(a, z) &:= \begin{cases} (g - c)z + (h - c)(a - z) & \text{if } z < a \\ (g - c)a & \text{if } z \geq a \end{cases} \\ &= \begin{cases} Gz - L(a - z) & \text{if } z < a \\ Ga & \text{if } z \geq a \end{cases} \\ &= G \min(z, a) - L(a - z)^+, \end{aligned} \quad (3.1)$$

where G and L are defined above. The representation (3.1) expresses the revenue in a single equation.

Consequently the expected revenue, denoted⁹ by $q(a)$, is given by

$$\begin{aligned} q(a) &:= E[r(a, Z)] = \sum_{z=0}^{\infty} r(a, z) f(z) \\ &= \sum_{z=0}^{a-1} (Gz - L(a - z)) P[Z = z] + Ga \sum_{z=a}^{\infty} P[Z = z]. \end{aligned} \quad (3.2)$$

The first term corresponds to the case when the demand z is less than the order quantity a . In this case, $a - z$ units of product remain unsold resulting in a loss of L units per item. The second term corresponds to the situation when demand is greater than or equal to the inventory. In this case, $z - a$ units of potential demand are lost and the reward is Ga .

Thus, the newsvendor problem can be summarized as that of finding

$$a^* \in \arg \max_{a \in A_0} q(a) = \arg \max_{a \in A_0} E[r(a, Z)],$$

where A_0 denotes the set of possible order quantities. The purpose of writing it this way is to prepare for its analysis using reinforcement learning methods in Chapters 10 and 11.

One particularly attractive feature of this model is that if the demand distribution is known, the optimal order quantity can be shown to satisfy

$$a^* = F^{-1} \left(\frac{G}{G + L} \right), \quad (3.3)$$

where $F(z) = \sum_{u=0}^z P[Z = u]$. Since the product is assumed to be available only in discrete units, a^* may need to be rounded up or down. Note that if demand is continuous this expression is exact.

⁹The quantity $q(a) := q(0, a)$ represents the state-action value function in the unique starting state 0.

The ratio $(G/G + L)$ is sometimes referred to as the *critical fractile*. It represents the fraction of demand that will be met by choosing a^* in this way. For example if $G = 2L$, two-thirds of the demand will be satisfied on average¹⁰.

Application challenges

The newsvendor model has been widely applied in the retail fashion-goods industry but also applies to other perishable goods. One novel application (unpublished) was its use by one us to the determine the number of instructors to have available for drop-in ski lessons each day. In applications, the greatest challenge is to model the demand distribution, which may vary by day-of-week or month-of-year and change over time.

3.2 Gridworld navigation

*A mathematician is a machine which turns coffee into theorems*¹¹.

Alfréd Rényi, Hungarian mathematician, 1921-1970.

As noted by Rényi, a working mathematician frequently requires coffee. To avoid time away from theorem proving, the mathematician may engage a robot to bring coffee when needed. In the example described here, the robot's task is to carry the mathematician's empty cup from the office to the coffee room, fill the cup with coffee and bring it back to the mathematician's office, all while avoiding barriers and pitfalls such as an open stairwell. When the mathematician is working on a particularly difficult theorem, a whole pot of coffee may be required, as shown by the robot on the book's cover and in Figure 3.3.

In general, a *Gridworld* represents a simplified environment commonly used to illustrate reinforcement learning methods. A Gridworld possesses a grid-like structure (like a chess board) where each cell in the grid represents a discrete state that an agent can occupy. The agent learns to move through the grid, typically aiming to achieve a goal while navigating around obstacles, collecting rewards, and avoiding penalties.

The model formulation in this section assumes that the robot *knows* the arrangement of the grid, which grid cell it occupies and the location of boundaries, target cells and penalties. This means the robot will never attempt to move outside the grid boundary. Alternative formulations assume:

1. the robot does not know its location,
2. the robot does not know the configuration of the grid, its boundaries, or the location of the stairwell,
3. the robot does not know the uncertainty involved when choosing which direction to move, or

¹⁰ Assuming the newsvendor faces this situation on many occasions.

¹¹ Wieschenberg [1999].

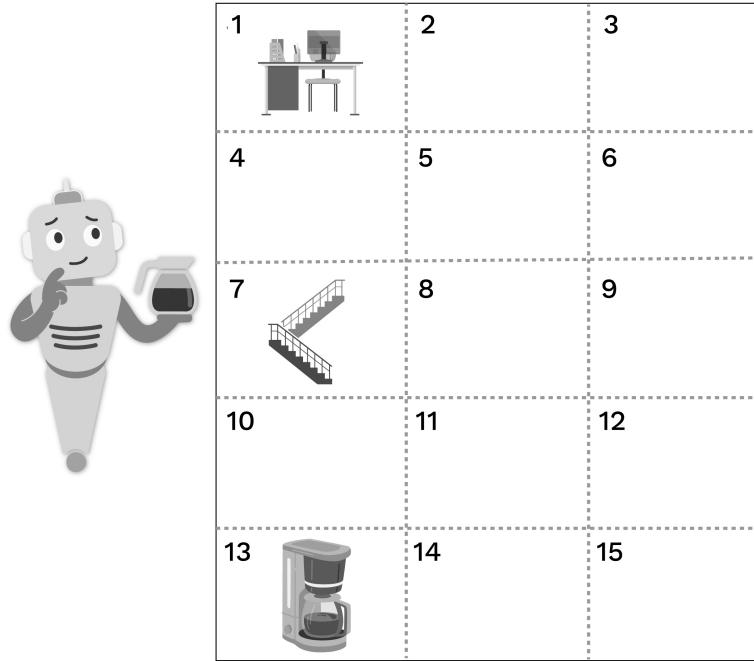


Figure 3.3: Schematic layout for Gridworld navigation example.

4. some combination of the above.

In these alternative situations, navigating through the grid becomes a *learning* problem. Such problems will be studied in depth in Chapters 10 and 11.

The problem description below assumes that the robot knows the grid configuration and which cell it occupies. In each cell except the stairwell, regardless of whether the coffee cup is empty or full, the robot can move in any of the four directions that does not take it outside the grid boundary. If the robot falls into the stairwell, or returns to the office with a full coffee cup, the episode terminates¹².

Assume movement on the grid is subject to uncertainty as follows. Let p_E and p_F denote the probability that the robot moves in its intended direction when the coffee cup is empty and full, respectively. If in some cell there are k possible locations where the robot can end up (including the current location), the probability the robot moves in each unintended direction or remains where it was is $(1 - p_E)/(k - 1)$ or $(1 - p_F)/(k - 1)$, depending on the status of the coffee cup. For example, suppose the robot has a full coffee cup and is in cell 6. Then if it intends to move down, it does so

¹²In some numerical examples in subsequent chapters, the episodes terminate only when the robot successfully delivers coffee.

with probability p_F and it moves left, moves up or remains in the same location with probability $(1 - p_F)/3$. The robot is more likely to be error prone with a full coffee cup because of the energy required not to spill the coffee. Thus, assume $p_E > p_F$.

The goal is to return with a full coffee cup as quickly as possible. If the robot successfully delivers coffee to the mathematician it receives a reward of B . If it falls down the stairs it incurs a penalty of X because the mathematician has to interrupt work to rescue the robot. Assume $X \gg B \gg 1$.

Decision epochs: Assume the process evolves in discrete time, where decision epochs correspond to the instant at which the robot decides in which direction to move. The first decision epoch after the robot enters the coffee room is used to fill up the cup, and the next one corresponds to the instant immediately after the cup has been filled and it decides where to move next. Thus

$$T = \{1, 2, \dots\}.$$

The set of decision epochs is unbounded because the robot continues attempting to deliver the coffee to the mathematician until the random time when it is either successful or falls down the stairs.

States: States represent the location of the robot in the numbered grid, plus an additional variable to indicate whether or not the coffee cup is empty (E) or full (F). The status of the coffee cup is required because it informs the success probability of an intended action. It also influences the direction in which the robot should proceed: a robot with an empty coffee cup seeks the coffee room, while a robot with a full coffee cup seeks the office. Therefore,

$$S = \{1, 2, \dots, 15\} \times \{E, F\}.$$

Actions: Actions represent the direction the robot attempts to travel. Assume the robot can only move up (north), down (south), right (east) and left (west), denoted by U, D, R and L , respectively. Under the assumption that the robot knows the layout of the grid and its location, the grid boundary constrains its intended movement so the set of permissible actions must take this into account. For example,

$$A_{(5,\cdot)} = \{U, D, R, L\}, \quad A_{(7,\cdot)} = \{U, D, R\}, \quad \text{and } A_{(15,\cdot)} = \{U, L\}.$$

Action sets for some particular states follow:

$$A_{(1,F)} = A_{(7,F)} = A_{(7,E)} = \{a_0\}, \quad A_{13,E} = \{a_1\}.$$

The states $(1, F)$, $(7, F)$ and $(7, E)$ are *absorbing* states. Once entered, the robot remains there forever; action a_0 corresponds to remaining in that state. Once such a state is entered, decision making stops. When the robot enters the coffee room with an empty cup, it takes one period to fill it. Denote the action of filling the cup by a_1 .

Rewards: Each intended movement action and the act of filling the coffee cup costs c . Often $c = 1$. Transitions into the office with a full coffee cup result in a reward of $B - c$. For example,

$$r((2, F), a, (1, F)) = B - c, \quad a \in A_{(2, F)} \quad (3.4)$$

Transitions into the stairwell receive a reward of $-X - c$ regardless of the status of the cup. For example,

$$r((4, k), a, (7, k)) = -X - c, \quad a \in A_{(4, k)} \text{ and } k \in \{E, F\}. \quad (3.5)$$

All other feasible state transitions between neighboring cells result in a reward of $-c$. For example,

$$r((5, k), a, (6, k)) = -c, \quad a \in A_{(5, k)} \text{ and } k \in \{E, F\}. \quad (3.6)$$

Finally, in this formulation, no rewards are received (or costs incurred) once the robot completes its task or falls down the stairs:

$$r((1, F), a_0, (1, F)) = 0 \text{ and } r((7, k), a_0, (7, k)) = 0, \quad k \in \{E, F\}. \quad (3.7)$$

Note that in the absence of uncertainty, it is easy to see that the robot can complete its task in 13 steps, so the maximum possible reward is $B - 13c$.

Transition probabilities: Some typical probabilities follow:

$$p((k, F)|(5, F), U) = \begin{cases} p_F & k = 2 \\ (1 - p_F)/4 & k \in \{4, 5, 6, 8\} \\ 0 & k \notin \{2, 4, 5, 6, 8\} \end{cases}$$

$$p((k, E)|(6, E), D) = \begin{cases} p_E & k = 9 \\ (1 - p_E)/3 & k \in \{5, 6, 9\} \\ 0 & k \notin \{3, 5, 6, 9\} \end{cases}$$

Transitions are deterministic when the system is in one of the absorbing states or when the robot enters the coffee room with an empty cup (since the only action is to fill the cup):

$$\begin{aligned} p((13, F)|(13, E), a_1) &= p((1, F)|(1, F), a_0) \\ &= p((7, F)|(7, F), a_0) = p((7, E)|(7, E), a_0) = 1. \end{aligned}$$

A key feature of this model is that the robot must trade off between safe but slow and risky but fast policies. This type of trade-off is characteristic of the key tension in many applications modeled by Markov decision processes. Here, by trying to reach the coffee room and returning to the office by the shortest route, the robot risks a high probability of falling down the stairs. If robot motion with an empty cup does not involve any randomness, that is $p_E = 1$, the robot will travel from the office to the coffee room by the shortest path but most likely take a more cautious path when returning to the office with a full cup.

Application challenges

The above example is artificial but Markov decision processes have been widely applied to robotic control. Realistic challenges include:

1. converting sensor readings to state variables,
2. transforming actions to motor commands,
3. representing variability in intended movement,
4. providing the robot with a mapping of the area, and developing reward structures.

These challenges are amplified when the application involves robotic movement in three-dimensional space.

3.3 Revenue management: Using price to manage demand

This section describes a different approach to inventory management. It pertains primarily to perishable products that decline in value over multiple periods.

As a concrete example, consider the challenges faced by a friend, Frank Z., previous owner of a chain of women's fashion stores in Vancouver, Canada. One evening, on the way to a poker session, he told Marty that:

Setting prices is like a game of chance; if I mark down prices too early in the season, I lose revenue, but if I wait too late, I can't sell my inventory and also must pay to store it.

A more general formulation of Frank's problem for a single product follows^[13]. A retailer acquires M units of a product (for example, a size 10 woman's dress in a particular pattern and design) and prices it at a_K . The retailer hopes that fashion-conscious and well-off customers will purchase all of the inventory at this price but if not, some inventory will remain. To sell the remaining inventory, the retailer may choose to *markdown* the price to make it accessible to more price-sensitive customers.

More formally, prices are set at the beginning of each month, are chosen from a finite set of prices, and remain constant throughout the month. Assume that when the price is a the demand in period n is random and Poisson distributed with rate $\lambda_{n,a}$. It is reasonable to assume that $\lambda_{n,a}$ decreases in both time (n) and price (a). The rationale is that fashion products become less attractive as time goes on so that expected demand at a fixed price decreases over time. Moreover, fundamental economic principles suggest that in any month demand should increase when price is reduced^[14].

^[13]Frank stocks hundreds of products each season so that this problem must be solved hundreds of times each season, a daunting task without a formal model to do so.

^[14]In economic modelling, this is referred to as a downward sloping demand curve.

Assume a monthly holding cost of $h(s)$ when the end of month inventory equals s units with $h(0) = 0$. Any goods left over at the end of month N are sold to an outlet store at a low price H per unit, representing the scrap value.

Decision epochs: Prices are set at the beginning of each month, so

$$T = \{1, 2, \dots, N\}.$$

States: States represent the number of items in stock at the start of each month in the planning horizon:

$$S = \{0, 1, \dots, M\}.$$

Actions: Actions represent the possible set of prices to choose from in each period. Assuming there are K candidate prices,

$$A_s = \{a_1, a_2, \dots, a_K\}$$

for $s \in \{1, \dots, M\}$. Note the price when $s = 0$ is irrelevant; set $A_0 = \{0\}$ for concreteness. Assume that $H < a_1 \leq a_2 \leq \dots \leq a_K$.

Rewards: If the inventory at the end of the month is j , that means $s - j \geq 0$ units were sold during that month. Thus, for $n < N$, $a_k \in A_s$, and $s \in S$,

$$r_n(s, a_k, j) = \begin{cases} a_k(s - j) - h(j), & j = 0, \dots, s \\ 0, & j = s + 1, \dots \end{cases}$$

and $r_N(s) = Hs$.

Transition probabilities: Since in this formulation no inventory is added during the planning horizon, only transitions to states with ending inventory $j \leq s$ have non-zero probability. If the demand equals or exceeds s , then the ending inventory is 0. This logic leads to the following transition probabilities for $s \in \{1, \dots, M\}$ and $n = 1, \dots, N - 1$:

$$p_n(j|s, a) = \begin{cases} e^{-\lambda_{n,a}} \lambda_{n,a}^{s-j} / (s - j)! & j = 1, \dots, s \\ \sum_{i=s}^{\infty} e^{-\lambda_{n,a}} \lambda_{n,a}^i / i! & j = 0 \\ 0 & j = s + 1, \dots \end{cases}$$

and¹⁵

$$p_n(j|0, 0) = \begin{cases} 1 & j = 0 \\ 0 & j \neq 0. \end{cases}$$

¹⁵Note the only action in state 0 is to set the price to 0.

Application challenges

Applying this model presents two challenges: determining the set of markdown prices and estimating the time-varying demand function parameters. In retail, the markdown prices can be set as a percentage of the original price such as “50% off” or “80% off”. Estimating the demand function requires considerable amounts of data and may be product specific. Data from similar products may provide guidance when there is insufficient historical data for the product. The parameter $\lambda_{n,a}$ may itself be represented as a function of n and a and learned during the decision problem.

3.4 Discrete-time queuing models

A queuing system consists of arrivals, a queue and one or more servers. Jobs arrive, wait in a queue if all servers are busy, are served by a free server and then depart the system when service is completed. Queuing systems have been well-studied in the operations research and engineering literature, and are applicable to a wide variety of service systems, including retail (jobs represent customers), healthcare (jobs represent patients), communication systems (jobs represent packets) and computer systems (jobs represent computing tasks).

From a decision perspective, the most widely studied models are:

1. **Service rate control:** The decision maker varies the service rate to control the queue length and throughput. The service rate may be controlled directly or through the addition and removal of servers.
2. **Admission control:** The decision maker chooses whether or not to admit an arriving job.
3. **Routing control:** In a network of queues, the decision maker chooses how to route jobs based on the workload at each queue.

The section formulates Markov decision process models for optimizing service rate and admission control. Exercise 8 provides a routing control example.

Some general comments regarding the formulation of discrete-time queuing systems follow:

1. Queuing systems are usually modeled as continuous-time Markov processes or semi-Markov processes. This example considers a discrete-time formulation. To do so, assume observation of the system starts at time 0. Let h denote a “small” unit of time and let decision epoch n correspond to “time” nh . Then the set of decision epochs denoted by $\{1, 2, \dots\}$ correspond to times $\{h, 2h, \dots\}$. The time increment h is chosen to be sufficiently small so that it is very unlikely more than one event (an arrival or service completion) occurs during a period of length h .

2. Queuing systems are usually modeled with infinite planning horizons to reflect that they are on-going and decision epochs occur frequently. No terminal reward is specified.
3. Rewards and transition probabilities are assumed to be independent of the decision epoch.
4. The system state is the number of jobs in the queue *and* in service.

3.4.1 Service rate control

Consider a single-server infinite capacity queuing system. At each decision epoch, the decision maker chooses a service completion probability from the set $\{a_1, a_2, \dots, a_K\}$. Let b denote the probability that a job arrives between two decision epochs, independent of the number of jobs in the system. Assume that $a_1 \leq a_2 \leq \dots \leq a_K$ and $a_K + b \leq 1$. Figure 3.4 provides a schematic representation of this queuing system.

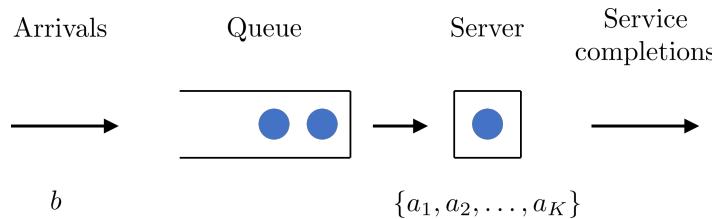


Figure 3.4: Schematic representation of a single server queuing system with adjustable service rate.

There are two costs to consider: a cost $m(a)$ per period for serving at rate a , and a delay cost of $f(s)$ per period when there are s jobs in the system at the start of the period. It makes sense to assume that both $m(a)$ and $f(s)$ are non-decreasing in their arguments.

A Markov decision process formulation follows. Note it assumes h is pre-specified and does not appear in the Markov decision process specification.

Decision epochs:

$$T = \{1, 2, \dots\},$$

corresponding to time points $\{h, 2h, \dots\}$.

States: States represent the number of jobs in the system (queue plus server):

$$S = \{0, 1, 2, \dots\}.$$

Actions: Actions represent the probability that a job currently being served is completed in the current period prior to the next decision epoch. For $s \in S$,

$$A_s = \{a_1, a_2, \dots, a_K\}^{16}.$$

Rewards: Costs should be represented as negative rewards, so

$$r(s, a) = -m(a) - f(s).$$

Note that this reward function is independent of the subsequent state.

Transition probabilities: For $s = 1, 2, \dots$ and $k = 1, 2, \dots, K$,

$$p(j|s, a_k) = \begin{cases} a_k & j = s - 1 \\ b & j = s + 1 \\ 1 - a_k - b & j = s. \end{cases} \quad (3.8)$$

For $s = 0$ and $k = 1, 2, \dots, K$,

$$p(j|0, a_k) = \begin{cases} b & j = 1 \\ 1 - b & j = 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3.9)$$

The rationale for these probabilities follows. When $s = 0$, a transition only occurs when there is an arrival. When $s > 0$, the assumption that h is small comes into play. This means that three things can happen, the state increases to $s + 1$ when an arrival occurs, the state decreases to $s - 1$ when a service is completed, and the state remains the same when neither occur.

A truncated model

An alternate formulation suitable for direct computation requires *truncating* the state space and modifying some transition probabilities. This is because direct numerical computation is impossible with a non-finite discrete state space¹⁷. In such a truncated model, all other model elements remain the same.

Truncated states: For some integer $W > 0$,

$$S = \{0, 1, \dots, W\}.$$

¹⁶Since the server is idle when there are no jobs in the system, it is optimal to use the slowest service rate so that one could set $A_0 = \{a_1\}$. Ideally the server should be turned off, but this formulation does not allow that. Moreover because h is small, it might be impractical to turn off the server.

¹⁷However, such computation may not be needed if the structure of an optimal policy and value function can be succinctly encoded, e.g., a control limit policy or monotonic value function. See Section 4.4.

Truncated transition probabilities: When $s = W$, the system becomes blocked so arrivals are not possible. This means that for $k = 1, \dots, K$,

$$p(s|W, a_k) = \begin{cases} a_k & \text{for } s = W - 1 \\ 1 - a_k & \text{for } s = W. \end{cases}$$

Choosing W presents some challenges. A well-controlled system will spend most of its time in low-occupancy states so an optimal policy might not be sensitive to the choice of W if it is chosen to be sufficiently large. A continuous approximation, where the mean and variance of the queue length can be computed in closed form, can provide some guidance. Alternatively, W may be varied and its impact noted.

3.4.2 Admission control

In admission control models, the decision maker assumes the role of a “gate keeper” by deciding whether or not to admit an arriving job into a queuing system. This provides an example of a model with a vector-valued state space and non-actionable states, which occur when no job arrives in the preceding period.

Again assume discretized time intervals of length h where h is sufficiently small so that at most one job can arrive between decision epochs and does so with probability b . If not admitted, the job is lost. Let $f(j)$ be the holding cost when there are j jobs in the system at the start of a period (after the admission decision, but before an arrival or service completion in the same period). Let w be the probability of a service completion between in a time interval of length h . Admitting a job into the system, generates a payment of R . Assume $b + w \leq 1$, which is reasonable when the time discretization step h is small. Figure 3.5 provides a schematic representation of the system and Figure 3.6 depicts the one-period dynamics.

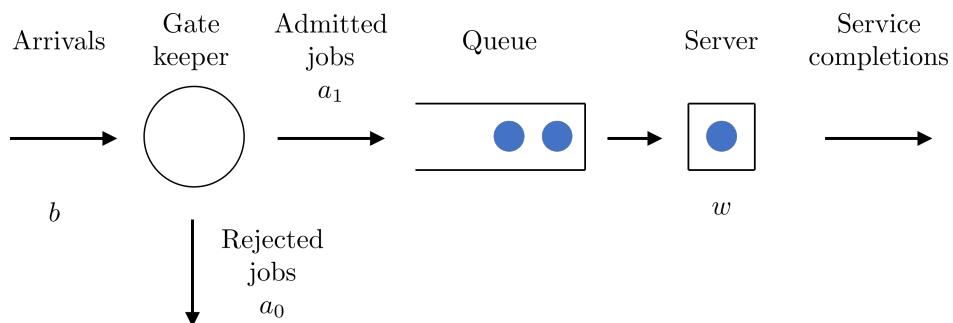


Figure 3.5: Schematic representation of a single server queuing system with admission control.

A Markov decision process formulation follows.

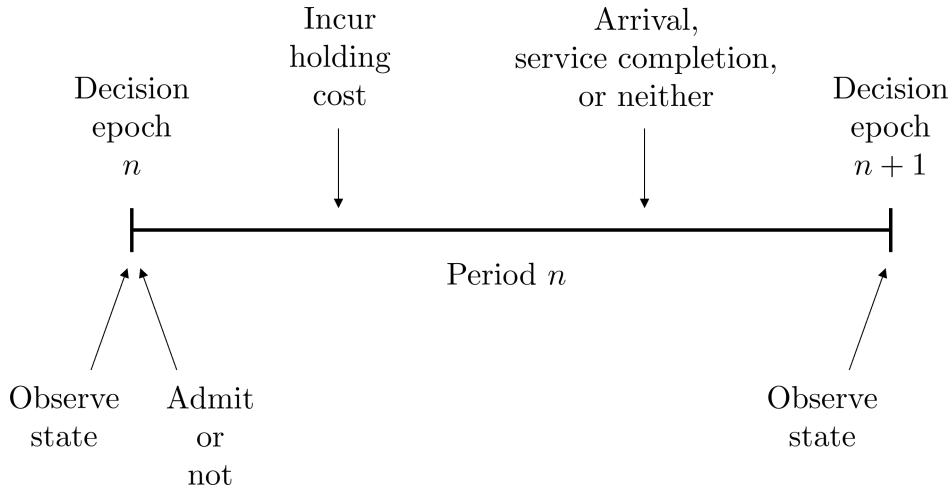


Figure 3.6: Timing of events in the queuing admission control model.

Decision epochs:

$$T = \{1, 2, \dots\}$$

corresponding to time points \$\{h, 2h, \dots\}\$.

States: The state has two components. The first component, denoted \$j\$, represents the number of jobs in the system and takes values in \$J = \{0, 1, \dots\}\$, and the second component, denoted \$k\$, indicates whether there is a job waiting for admission (\$k = 1\$) or not (\$k = 0\$) so that \$k \in \{0, 1\}\$. Then the state space can be written as

$$S = J \times \{0, 1\},$$

with a typical state represented by the two-dimensional vector \$(j, k)\$.

Actions: Let \$a_0\$ correspond to the action “do not admit” and \$a_1\$ to the action “admit”¹⁸. Since admission is possible only if an arrival occurred since the previous decision epoch, there is no choice in states¹⁹ where \$k = 0\$. Thus, for \$j = 0, 1, \dots\$,

$$A_{(j,k)} = \begin{cases} \{a_0, a_1\}, & k = 1 \\ \{a_0\}, & k = 0 \end{cases}$$

Recall that to formulate a Markov decision process model, actions need to be specified in all states even when the action set contains a single element.

¹⁸As a general convention, when the action space comprises two actions where one action is to “do nothing”, that action will be referred to as \$a_0\$ and the other as \$a_1\$.

¹⁹Such states are said to be *non-actionable*.

Rewards: For $j = 0, 1, \dots$,

$$r((j, k), a) = \begin{cases} R - f(j+1), & k = 1, a = a_1 \\ -f(j), & k = 0, a = a_0. \end{cases}$$

Note that in this model, the rewards are independent of the subsequent state (j', k') .

Transition probabilities: If the action is to not admit ($a = a_0$), then whether there is a job currently waiting to be admitted or not ($k = 0$ or 1) is irrelevant since a waiting job will not be admitted. Thus, when $a = a_0$ for $j = 1, 2, \dots$ and $k = 0$ or 1 ,

$$p((j', k')|(j, k), a_0) = \begin{cases} w & j' = j-1, k' = 0 \\ b & j' = j, k' = 1 \\ 1-b-w & j' = j, k' = 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3.10)$$

Since service completions are not possible if the system is empty, the transition probabilities when $j = 0$ and $a = a_0$ are given by

$$p((j', k')|(0, k), a_0) = \begin{cases} b & j' = 0, k' = 1, \\ 1-b & j' = 0, k' = 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3.11)$$

The “admit” action a_1 applies only when $k = 1$. So for $j = 0, 1, 2, \dots$

$$p((j', k')|(j, k), a_1) = \begin{cases} w & j' = j, k' = 0 \\ b & j' = j+1, k' = 1 \\ 1-b-w & j' = j+1, k' = 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3.12)$$

In contrast to (3.10), (3.12) is valid when $j = 0$ since even when the system is empty at the decision epoch, a decision to admit will add a job to the system, which can be completed during the same period with probability w .

Figure 3.7 summarizes the possible state transitions for each action. Despite the simplifying assumption that at most one event can happen in a period, keeping track of all transitions requires a careful accounting of events and actions and their interactions. This example is revisited in Chapter 4, which provides a simpler formulation of the model in terms of the *post-decision state* and *state-action value functions*. The post-decision state provides an alternative view of the decision timeline depicted in Figure 3.6 that allows the transition dynamics to be modeled more easily. As emphasized frequently in this chapter, drawing a correct timeline is an important part of formulating the model correctly.

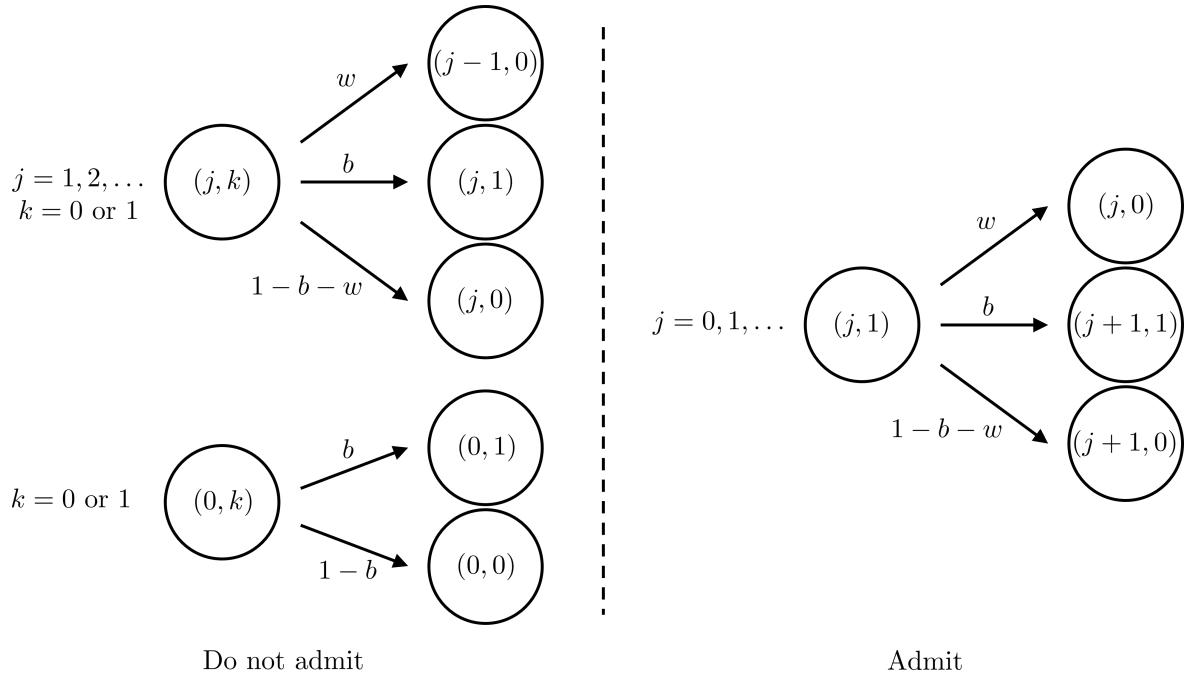


Figure 3.7: Possible state transitions for admission control problem. Note that not admitting a waiting job results in the same transitions as when no job arrives.

A post-decision state formulation

The following formulation is based on *post-decision states* in contrast to the above formulation based on *pre-decision states*.

Shifting the timing of decision epochs leads to a simpler formulation. Suppose instead that decision epochs occur just prior to resolving the uncertain event “Arrival, service completion or neither” in Figure 3.6 so the timing now follows Figure 3.8. In this formulation at a decision epoch, the previous decision has been implemented so that the state need only represent the number of jobs in the system at this decision epoch. Actions need to be modified to be contingent on the resolution of the uncertain event as described below. Note that the holding cost is assessed after the action has been implemented.

Decision epochs:

$$T = \{1, 2, \dots\}$$

corresponding to time points $\{h, 2h, \dots\}$, where $t = 1$ represents the time of the first decision.

States:

$$S = \{0, 1, \dots\}$$

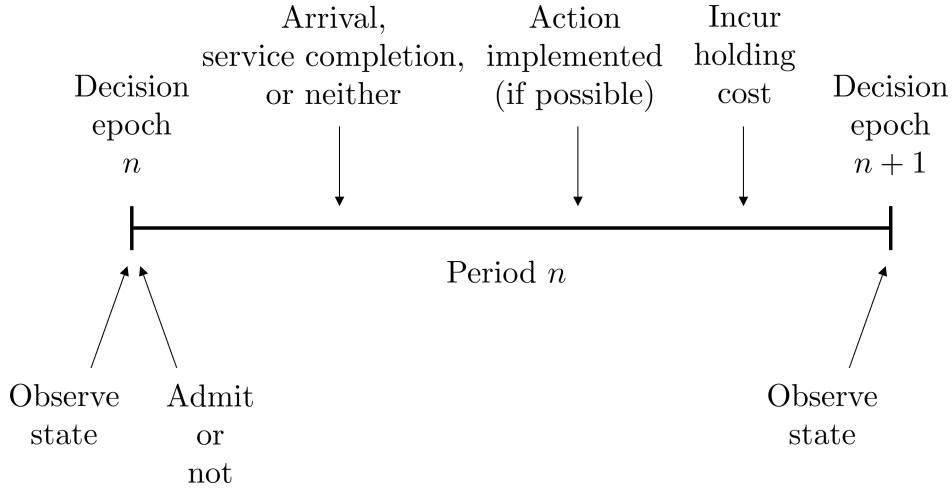


Figure 3.8: Timing of events in alternative formulation of the queuing admission control model. The expression “action implemented (if possible)” applies only if there is a prior arrival. Otherwise it cannot be implemented.

where $s \in S$ represents the number of jobs in the system at a decision epoch.

Actions: Let a_0 correspond to the action “do not admit” and a_1 to the action “admit”. Since the decision maker does not know whether or not there will be an arrival, the result of this action must be contingent on the resolution of the event. Thus a_0 must be interpreted as “do not admit if there is an arrival” and a_1 as “admit if there is an arrival”. By default if there is no arrival, the action a_0 applies since it will have the same effect as when there is an arrival. Thus for all $s \in S$,

$$A_s = \{a_0, a_1\}.$$

Rewards: In this formulation, rewards depend on the subsequent state and when the holding cost is assessed. Assuming the holding cost is incurred after the uncertain event is resolved and the decision implemented yields:

$$r(s, a_0, j) = \begin{cases} -f(j) & \text{for } j = s, s-1 \text{ and } s > 0 \\ -f(0) & \text{for } j = 0 \text{ and } s = 0 \\ 0 & \text{otherwise} \end{cases}$$

$$r(s, a_1, j) = \begin{cases} R - f(j+1) & \text{for } j = s+1 \text{ and } s \geq 0 \\ -f(j) & \text{for } j = s, s-1 \text{ and } s > 0 \\ -f(0) & \text{for } j = 0 \text{ and } s = 0 \\ 0 & \text{otherwise} \end{cases} \quad (3.13)$$

Observe that the payment R is received only when choosing a_1 and an arrival occurs, which corresponds to the case when $j = s + 1$ in (3.13). In all other cases there are no admissions and only the holding cost is incurred.

Transition probabilities: Under a_0 a service completion is the only possible event that can occur when $s > 0$. It results in a transition to $s - 1$. When the system is empty, such a transition cannot occur so that:

$$p(j|s, a_0) = \begin{cases} 1-w & \text{for } j = s \text{ and } s > 0 \\ w & \text{for } j = s - 1 \text{ and } s > 0 \\ 1 & \text{for } j = 0 \text{ and } s = 0 \\ 0 & \text{otherwise} \end{cases}$$

Under a_1 , an arrival results in a transition to $s + 1$. When there is no arrival the system either remains in the same state or moves to state $s - 1$ if $s > 0$ if there is a service completion.

$$p(j|s, a_1) = \begin{cases} b & \text{for } j = s + 1 \text{ and } s \geq 0 \\ w & \text{for } j = s - 1 \text{ and } s > 0 \\ 1 - b - w & \text{for } j = s \text{ and } s > 0 \\ 1 - b & \text{for } j = 0 \text{ and } s = 0 \\ 0 & \text{otherwise} \end{cases} \quad (3.14)$$

Note that this formulation may be more natural when using state-action value functions and simulation or when formulating it directly as a continuous time model. It is left to the reader to see which is preferred from the perspectives of ease of formulation and computation.

Application challenges

Applying these models requires estimates of arrival probabilities, service probabilities and costs. Queuing models are more commonly formulated in continuous time with inter-arrival and service times modeled using exponential random variables. Thus, if arrivals occur at rate λ , the probability of one arrival in an interval of length h is given by $\lambda h + o(h)$, the probability of no arrivals in an interval of length h is $1 - \lambda h + o(h)$, and the probability of greater than one arrival in an interval of length h is $o(h)$ where $o(h)$ is an expression that converges to zero as h decreases to zero. Thus, it is convenient to set the probability of an arrival in a short time interval of length h to be λh .

As in other models, determining costs and rewards is somewhat arbitrary, and may depend heavily on the application. It is important to investigate the impact of specific choices through sensitivity analyses.

3.5 Behavioral decision making: When should a lion hunt?

Markov decision processes provide a natural framework for modeling behavior when an organism faces a decision that trades off survival with reserving energy. Examples include choosing a location for food acquisition, deciding when to hunt for food, choosing a group size when hunting and deciding when to abandon its offspring. The primary objective in such research is the determine whether results from an optimization model agree with observed animal behavior.

As an example, consider the challenge facing a lion (*panthera leo*) when deciding to hunt for food. Suppose the lion seeks to maximize its probability of survival over a season of N days. A mature lion has an energy storage capacity of C units. Each day it does not hunt the lion depletes its energy reserves by d units. Hunting requires h units of energy with $h > d$. If its energy reserves fall below c_0 units, it will not survive to the next day.

At the start of each day, the lion decides whether or not to hunt and if so, what prey to seek. Typically, lions hunt impalas, gazelles, wildebeests, giraffes and zebras. About half of the time they hunt in groups. This development assumes that the lion hunts alone²⁰. Catch probability varies with species hunted. Assume that there are M species to choose from. Let w_m denote the probability that the lion catches an animal of species m , $m = 1, 2, \dots, M$. A successful hunt for species m yields e_m units of energy. For simplicity, assume all relevant quantities have been discretized.

Decision epochs: Decisions are made at the start of each day during the season, so

$$T = \{1, 2, \dots, N\}.$$

States: The state represents the lion's energy reserves at each decision epoch. Thus

$$S = \{0, 1, \dots, C\}.$$

Actions: Actions in state s may be denoted by

$$A_s = \begin{cases} \{a_0, a_1, \dots, a_M\} & s \in \{c_0, c_0 + 1, \dots, C\} \\ \{a_0\} & s \in \{0, 1, \dots, c_0 - 1\}, \end{cases}$$

where action a_0 corresponds to “do not hunt” and action a_m corresponds to “hunt for species m ”.

²⁰Exercise 19 asks you to formulate the group size decision problem.

Rewards: Given the lion's survival objective, the lion receives a reward of 1 if alive at the end of the season and a reward 0 if not. Therefore

$$r_N(s) = \begin{cases} 1 & s \in \{c_0, c_0 + 1, \dots, C\} \\ 0 & s \in \{0, 1, \dots, c_0 - 1\}. \end{cases}$$

No rewards are accrued throughout the planning horizon, so $r_n(s, a, j) = 0$ for $n = 1, 2, \dots, N - 1$, $a \in A_s$, $s \in S$ and $j \in S$.

Transition probabilities: When the lion has energy reserves of s units at the start of the day, hunts for species m and is successful, its energy reserves at the start of the next day is $\min\{s - h + e_m, C\}$. The logic underlying this observation is that starting in state s , the lion uses h units of energy hunting and if successful obtains prey that provides e_m units of energy. However, the total energy acquired cannot exceed its capacity C . As the result of an unsuccessful hunt, its energy reserves fall to $\max\{s - h, 0\}$. Therefore

$$p_n(j|s, a) = \begin{cases} 1 & j = \max\{s - d, 0\}, s = c_0, \dots, C, a = a_0 \\ w_m & j = \min\{s - h + e_m, C\}, s = c_0, \dots, C, a = a_m \\ 1 - w_m & j = \max\{s - h, 0\}, s = c_0, \dots, C, a = a_m \\ 1 & j = s, s = 0, 1, \dots, c_0 - 1, a = a_0 \\ 0 & \text{otherwise.} \end{cases}$$

Note that the transition probabilities take into account the fact that the lion's energy level cannot exceed the maximum capacity or fall below 0, and if its energy level falls below c_0 , it cannot hunt.

Application challenges

This application highlights the fact that it can be challenging to determine parameter values for a Markov decision process, and to do so, one must often appeal to a wide range of sources. Moreover, in this particular example, it is essential to understand the underlying animal behavior and model it correctly, ideally with expert input. An added benefit of developing a formal Markov decision process model is that it identifies relevant parameters that can motivate related research.

The ecology literature suggests values for many of the key model parameters although not always exactly in the form needed. One can use $C = 30$ and $d = 6$ kilograms, based on averages of male and female lions. If a lion-specific parameter is not available from the literature, borrowing values from other species may provide some guidance. For example, wild dogs expend about 23% more energy per hour when hunting, with the average duration of a hunt approximately equal to 40 minutes. Noting that lions undertake on average three chases per day and assuming that they expend

the same incremental amount of energy while hunting, on a day they decide to hunt, they will spend 3% more energy than on day they decide to rest so that $h = 1.03d$.

The literature suggests that gazelles yield a mean biomass of 12 kilograms with a catch probability of 0.15 on a single hunt. Observational data suggests that a lion may hunt up to three times a day if earlier hunts are unsuccessful; such dynamics can be incorporated into our model. Modeling the hunting of zebras presents additional challenges. Zebras yield an estimated 164 kilograms of edible biomass with a catch probability between 0.15 and 0.19 depending on the hunt location. Since they are large, their carcasses last for several days and are shared among several lions. Determining how much is available for the hunter requires further assumptions. Other sources provide estimates of the edible biomass for other types of prey such as impalas (29 kg), wildebeests (150 kg), and giraffes (468 kg), but not catch probabilities, which are harder to estimate. It is quite common for domain-specific literature to report data that allows us to estimate only some of the parameters of a Markov decision process, since such data is typically reported without a decision-making objective in mind. As a result, many assumptions must often be made, as described above. Especially when estimates are quite variable and many assumptions are needed, it is recommended to conduct sensitivity analyses of these parameters.

3.6 Clinical decision making: An application to liver transplantation

Markov decision processes have been widely applied to medical decision problems including organ transplantation, HIV treatment, cholesterol management and cancer diagnostics. As an illustration, the following describes a decision problem that arises when a patient requires a liver transplant.

This section describes the decision process for a patient with end-stage liver disease (ESLD)²¹ who requires a liver transplant. Organs intermittently become available²² and vary in quality. Depending upon the quality of the organ, the patient's medical team may either accept the organ and transplant it immediately, or reject it and wait for a higher quality organ.

To make this concrete, a relatively healthy patient may reject a lower-quality organ in the hopes of being offered a higher-quality organ in the future. On the other hand, a patient in poor health may accept the first liver available. A Markov decision process model can be used to formalize this decision problem and explore trade-offs.

To facilitate modeling, patient health status and organ quality are represented by discrete categories. Health states vary from 1 to H where 1 represents the healthiest state and H represents the least healthy state. The state Δ represents death by liver

²¹*End-stage liver disease* or *cirrhosis* refers to a condition in which a patient's liver is severely damaged and no longer able to function adequately. Without a transplant, it is usually fatal.

²²From a recently deceased individual. These are referred to as *cadaveric* organs in the medical literature.

failure. Similarly, liver quality states are ordered from 1 (highest) to L (lowest). The state Φ represents an occasion when no liver is available.

Rewards measure life expectancy in days. On a day when there is no transplant, the patient accrues a reward of 1 (an extra day of life) independent of the health state. The (discounted) life expectancy post-transplant of a patient in health state h who accepts a liver of quality l is represented by $R(h, l)$. It is reasonable to assume (why?) that $R(h, l)$ is non-increasing in h and l .

To facilitate modeling, assume decisions are made daily, immediately after determining whether a liver is available and, if so, its quality. If a patient does not receive a transplant, their health state either remains stable or worsens. Furthermore, the likelihood of a patient being offered a liver depends on their current health state.

At each decision epoch, let

- $\gamma(h'|h)$ denote the probability that a patient in health state h deteriorates to state h' with $h' = h, h + 1, \dots, H$ at the next decision epoch,
- $\delta(h)$ denote the probability a patient in health state h dies from liver failure prior to the next decision epoch,
- $\beta(l|h)$ denote the probability that a patient in health state h is offered a liver of quality $l = 1, \dots, L$ immediately prior to the next decision epoch, and
- $\phi(h)$ denote the probability a patient in health state h is *not* offered a liver prior to the next decision epoch.

Assume these distributions are stationary and independent. A Markov decision process formulation follows.

Decision epochs: Assume decisions are made daily after ascertaining whether a liver has been offered and if so, its quality. Further assume an infinite²³ horizon. Thus,

$$T = \{1, 2, \dots\}.$$

States: States represent the patient's health (if alive) and liver quality (if available). Let the absorbing state Γ represent the post-transplant state. For convenience, define $S_H = \{1, \dots, H\}$, $S_L = \{1, \dots, L\}$ and $S_L^+ = S_L \cup \{\Phi\}$. Then,

$$S = (S_H \times S_L^+) \cup \{\Delta, \Gamma\}.$$

Note there is no need to distinguish Δ and Γ since the decision process ends in either case.

²³A practical upper bound on the number of decision epochs might be 26,000 (assuming no transplants for people over 90 or younger than 20). Given this upper bound, an infinite horizon model may be appropriate and simpler, especially since the process will reach an absorbing state eventually, either post-transplant or death most likely before reaching 90 years old.

Some explanation may be helpful. States of the form $(h, l) \in S_H \times S_L$, correspond to being in health state h and having an organ of quality l available. Φ replaces L when no organ is available and $s = \Gamma$ when no more decisions are possible.

Actions: Let a_1 represent the action to accept an organ (assuming one is available) and a_0 represent the action to wait (do not accept an organ). Also let a_0 represent the “do nothing” action, which applies in the death state, in the post-transplant state, and in any state when no organ is offered. Hence

$$A_s = \begin{cases} \{a_0, a_1\} & s \in S_H \times S_L \\ \{a_0\} & s \in (S_H \times \{\Phi\}) \cup \{\Delta, \Gamma\}. \end{cases}$$

Note that a_1 is applicable only when the patient is alive and an organ is available, whereas a_0 applies regardless of whether an organ is available, the patient has passed away, or a transplant has already occurred.

Rewards: The reward function may be represented for $s = (h, l)$ as:

$$r(s, a, j) = \begin{cases} R(h, l) & (h, l) \in S_H \times S_L, a = a_1, j = \Gamma \\ 1 & (h, l) \in S_H \times S_L^+, a = a_0, j = (h', l') \in \{h, \dots, H\} \times S_L^+ \\ 0 & (h, l) \in S_H \times S_L^+, a = a_0, j = (h', l') \in \{1, \dots, h-1\} \times S_L^+ \\ 0 & (h, l) \in S_H \times S_L^+, a = a_0, j = \Delta \\ 0 & s \in \{\Delta, \Gamma\}, a = a_0, j = s. \end{cases}$$

Note the second and third expressions correspond to the assumption that the patient’s health state cannot improve²⁴. Further the fourth expression indicates that the patient does not accrue an additional day of life if death occurs prior to the next decision epoch.

Transition probabilities: If the patient does not receive a transplant, then a transition from health state h to state h' with $h' \geq h$ may occur²⁵. If the patient chooses to receive a transplant, then there is a deterministic transition to the post-transplant state, Γ . In the post-transplant or death state, the system remains in that state forever. Therefore

²⁴Note that the reward under a transition to a $j = (h', l') \in \{1, \dots, h-1\} \times S_L^+$ is arbitrary since this transition cannot occur.

²⁵Recall that smaller values of h correspond to better health.

$$p(j|s, a) = \begin{cases} \gamma(h'|h)\beta(l'|h) & s = (h, l) \in S_H \times S_L^+, a = a_0, j = (h', l') \in \{h, \dots, H\} \times S_L \\ \gamma(h'|h)\phi(h) & s = (h, l) \in S_H \times S_L^+, a = a_0, j = (h', l') \in \{h, \dots, H\} \times \Phi \\ \delta(h) & s = (h, l) \in S_H \times S_L^+, a = a_0, j = \Delta \\ 1 & s \in S_H \times S_L, a = a_1, j = \Gamma \\ 1 & s \in \{\Delta, \Gamma\}, a = a_0, j = s \\ 0 & \text{otherwise.} \end{cases}$$

Application challenges

Application of Markov decision process models to clinical decision making requires detailed medical domain knowledge including the nature and progression of the disease, and the treatment options and processes. For instance, applying this model to liver transplantation requires in-depth knowledge of the liver transplant system, the process used to allocate organs and the progression of end-stage liver disease. Data required includes discretized patient health and liver quality states, an estimate of post-transplant life expectancy, probabilities of health state deterioration, and arrival distributions of organs for transplantation by quality. Such data may be obtained from transplant centers and organizations that manage the transplantation system, such as the United Network for Organ Sharing (UNOS) in the United States. Patient health status can be measured using the Model for End Stage Liver Disease (MELD) score, which is a function of various laboratory values. Higher scores indicate poorer health and a higher mortality rate²⁶. When data is sparse, MELD scores can be aggregated or smoothed. A similar approach can be taken when defining liver quality states, which may depend on donor age, race and sex. UNOS data may be used to estimate $R(h, l)$ (days of survival post-transplant) with a proportional hazards model²⁷ and organ arrival rates. Transitions between health states can be modeled using a natural history model²⁸.

3.7 Advance appointment scheduling

In many applications, decision makers must allocate scarce resources prior to the arrival of future random demand. The formulation below is from the perspective of a scheduler in a hospital diagnostic imaging department who faces the challenge of scheduling appointments for current medical imaging requests to meet patient-specific clinical wait time targets, without knowing how many and when future requests will arrive. Ideally,

²⁶Internet sources suggest that the 3-month survival rate of 27% in patients with the highest MELD score of 40 and 98% for patients with MELD score between 1 and 9.

²⁷This is a statistical model that can be used to determine the impact of covariates on survival times when some patients in the data set are still alive.

²⁸A commonly used epidemiological model that simulates disease progression accounting for different risk factors.

the scheduler would strive to minimize the fraction of patients scheduled beyond their target dates. The formulation²⁹ below introduces a cost to achieve this objective.

Suppose appointment requests arrive throughout the day and at the end of each day a radiologist assigns each request to one of K *urgency classes*. Urgency class k is associated with a target wait time of T_k days, $k = 1, 2, \dots, K$, which is chosen based on clinical considerations. A patient in urgency class k should be scheduled prior to day T_k . The urgency classes are ordered, with 1 having the highest priority, so $T_1 < T_2 < \dots < T_K$. If a patient in urgency class k is scheduled after T_k , a cost³⁰ C_k is incurred proportional to the number of days past T_k the appointment is scheduled. This cost can be thought of as a penalty related to worse clinical outcomes due to delayed diagnosis. Assuming $C_1 > C_2 > \dots > C_K$ corresponds to setting higher delay costs for the most urgent cases. If a patient in class k is scheduled before the target T_k , no cost is incurred.

Let $p_k(w)$, $k = 1, 2, \dots, K$ denote the probability that w new class k appointments arrive each day. Let $\mathcal{W} = \{0, 1, \dots, M\}$ denote the set of possible values of w . To ensure a finite state and action formulation, assume M is finite. Daily capacity is divided into B appointment slots. This means that at most B regular time appointments can be booked each day. Assume there are an unlimited number of overtime slots available and the system incurs a cost of h for each patient scheduled to overtime. Implicit is that $h > C_1$, which means that delaying a patient by a day beyond the target time is less costly than scheduling the patient to overtime. But for a sufficiently large delay, the cost of overtime will be less than the cost of delay.

Once all requests have an assigned urgency class, the scheduler assigns an appointment date to each request and informs the patient. Figure 3.9 provides a timeline for this process. The challenge is to schedule today's requests before realizing future requests in the face of limited capacity.

There are numerous issues arising in this and other similar scheduling problems that can impact modeling:

1. Does the system have access to surge capacity or overtime?
2. How does appointment length vary between patients?
3. Can overbooking be used to account for patient no-shows and late cancellations?
4. Can appointment dates be changed after scheduling?
5. Are the targets flexible or must they be met?
6. Does demand have any seasonal patterns or correlation across urgency classes?

²⁹Such an application presents many modeling challenges and requires formulating the problem carefully. The formulation described below was not immediately obvious to the investigators (see Section 3.10) and required considerable trial and error to achieve.

³⁰Note these costs are artificial, wait time targets in medical settings are often specified in terms of the proportion of demand that is scheduled prior to its target date.

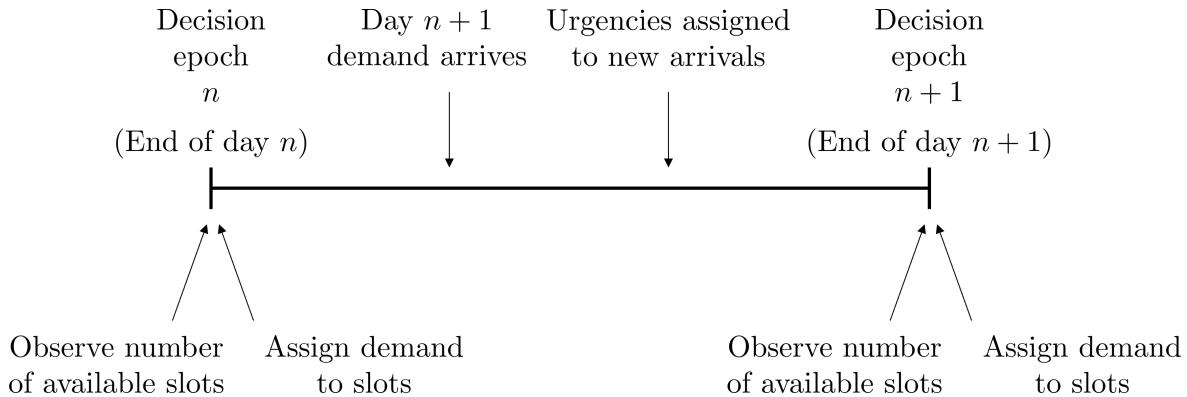


Figure 3.9: Timing of events in the advance appointment scheduling model.

7. How far in advance can appointments be scheduled?

The formulation below assumes:

- access to overtime,
- fixed appointment lengths,
- no cancellations,
- no rescheduling,
- flexible target dates (but with a penalty for exceeding the target date),
- stationary arrivals and uncorrelated demand between urgency classes,
- a fixed *booking horizon* of N days.

Note that the booking horizon refers to how far into the future current appointment requests can be scheduled, and not the length of the planning horizon. As an alternative to overtime, appointments not scheduled during a particular day may be held over for scheduling in the future at some cost. Exercise 16 considers this variation.

Decision epochs: Decision epochs correspond to the time in the day when the scheduling clerk assigns an appointment date to each appointment request waiting to be scheduled. This is naturally modeled as an infinite horizon problem, so

$$T = \{1, 2, \dots\}.$$

States: A typical state of the system is represented by $s = (b_1, \dots, b_N, w_1, \dots, w_K)$ where $b_n \in \mathcal{B} = \{0, 1, \dots, B\}$ denotes the number of appointments that have already been booked on day n for $n = 1, \dots, N$ and $w_k \in \mathcal{W}$ denotes the number of appointments of urgency class k waiting to be scheduled at a decision epoch. Note that if there are b_n appointments booked on day n , then there are $B - b_n$ remaining appointment slots on that day.

To simplify notation, introduce the vectors $\mathbf{b} := (b_1, \dots, b_N)$ and $\mathbf{w} := (w_1, \dots, w_K)$. Thus, a state is denoted by (\mathbf{b}, \mathbf{w}) , and the state space is

$$S = \mathcal{B}^N \times \mathcal{W}^K.$$

Actions: Actions represent the number of waiting patients in urgency class k to schedule on each day within the booking window and possibly through overtime if B appointments if advantageous. Let x_{kn} denote the number of class k patients to book on day n and y_k denote the number of class k patients to book for overtime the next day. Note that it is not necessary to consider booking overtime slots at some point in the future since there is no limit on the number of patients booked through overtime, and booking overtime further in the future would simply increase costs.

Define the vectors $\mathbf{x} := (x_{11}, \dots, x_{1N}, x_{21}, \dots, x_{2N}, \dots, x_{K1}, \dots, x_{KN})$, $\mathbf{y} := (y_1, \dots, y_K)$ and $\mathbf{0} := (0, 0, \dots, 0)$ with lengths NK , K and $(N+1)K$, respectively. The action set $A_{(\mathbf{b}, \mathbf{w})}$ is

$$A_{(\mathbf{b}, \mathbf{w})} = \left\{ (\mathbf{x}, \mathbf{y}) \geq \mathbf{0} \left| \begin{array}{l} \sum_{n=1}^N x_{kn} + y_k = w_k \text{ for } k = 1, \dots, K \text{ and} \\ b_n + \sum_{k=1}^K x_{kn} \leq B \text{ for } n = 1, \dots, N \end{array} \right. \right\}. \quad (3.15)$$

The first condition in (3.15) ensures that all class k requests must be scheduled either to a specific day or to overtime. The second condition ensures that at most B patients may be scheduled to regular time each day.

Rewards: Assume that the penalty associated with scheduling an urgency class k request n days from the current day costs $C_k(n - T_k)^+$. This function specifically models the cost as linear in the number of days a class k appointment is scheduled beyond its target, while incurring zero costs for scheduling prior to the target. Exercise 17 considers the variation where instead of the target representing a fixed day, a target window with both and earliest and latest time is assumed.³¹

³¹This applies to chemotherapy schedules where there is a narrow window when a treatment can be given.

The reward for choosing actions (\mathbf{x}, \mathbf{y}) in state (\mathbf{b}, \mathbf{w}) is

$$r((\mathbf{b}, \mathbf{w}), (\mathbf{x}, \mathbf{y})) = - \sum_{k=1}^K \sum_{n=1}^N C_k (n - T_k)^+ x_{kn} - h \sum_{k=1}^K y_k.$$

The reward (negative of cost) captures the costs associated with exceeding the target times as well as overtime costs for the current set of appointment requests. Observe that the reward does not depend on the next state since the subsequent \mathbf{b} vector is a deterministic function of the current action.

Transition probabilities: Transition probabilities depend on both action choice and the random arrival distribution. At the start of a period, the calendar rolls forward one day so previous bookings that were n days from the previous decision epoch are now $n - 1$ days from the current decision epoch. Added to these bookings are the newly arriving demand that is booked over the N -day horizon starting at the current decision epoch. Finally, since there were no bookings $N + 1$ days from the previous decision epoch, there are 0 appointments booked N days from the current decision epoch. The demand that has arrived since the last decision epoch is the only stochastic element in this model. Once this random demand has been assigned to urgency classes, the probability transitions are:

$$p((\mathbf{b}', \mathbf{w}') \mid (\mathbf{b}, \mathbf{w}), (\mathbf{x}, \mathbf{y})) = \begin{cases} \prod_{k=1}^K p_k(w'_k) & \mathbf{b}' = (b_2 + \sum_{k=1}^K x_{k2}, b_3 + \sum_{k=1}^K x_{k3}, \dots, b_N + \sum_{k=1}^K x_{kN}, 0) \\ 0 & \text{otherwise.} \end{cases} \quad (3.16)$$

Some comments on this formulation follow:

1. The concept of a booking horizon may be regarded as an artifact of the modeling process and imposed to maintain a finite state space. Since overtime is unlimited, a fixed booking horizon applies. Without overtime, an unbounded booking horizon may be needed.
2. Because the booking horizon remains constant between decision epochs, but moves forward each day, the model may be regarded as using a *rolling horizon* model.
3. Note that the transitions decompose into a deterministic part corresponding to the number of booked appointments each day and a random part corresponding to the random demand for each urgency class.
4. When the maximum daily demand for each urgency class is M , the model has $(C + 1)^N(M + 1)^K$ states. This makes direct computation infeasible for practical sizes of these parameters and motivates the need for approximation (see Chapter 9).

5. After an action is implemented at decision epoch n , the \mathbf{b} component of the state does not change until after decision epoch $n+1$. For reasons discussed in Section 4.5 and Chapter 11, it may be more convenient to formulate the model in terms of post-decision states.
6. In reality, many possible reward structures may be applicable for this problem. For example, as an alternative to incurring costs for appointments scheduled beyond their target date, a decision maker could strive to maximize the fraction of patients scheduled within their target dates.

Application challenges

Application challenges include specifying a booking horizon, specifying how unscheduled demand is satisfied, determining urgency classes and targets, specifying costs for delays, and estimating demand.

In real problem settings, a booking horizon may be determined by the decision maker based on their typical clinical processes. It has also been shown that when unlimited appointment diversion is possible, such as through overtime, an optimal policy is independent of the booking horizon, provided it exceeds the largest wait time target. Urgency classes should be defined based on clinical guidelines by medical professionals. The above model formulation was based on a real application. In that application the classes were “urgent” (7 day target), “semi-urgent” (14 day target) and “non-urgent” (28 day target). Emergency cases were scheduled to a different resource. Relative delay costs can potentially be quantified by calculating the impact on clinical outcomes of delayed treatment due to the delayed imaging.

Future demand may be forecasted using historical data. In practice, these distributions may be non-stationary as volumes generally increase over time. As an alternative to estimating the demand for each urgency class separately, if the historical data suggests that the relative proportion of cases from each urgency class is stable, it may be best to estimate total demand and then split it among each class based on the fixed proportion.

3.8 Optimal stopping

An elegant collection of applications are referred to as “optimal stopping problems”. Examples include selling an asset, finding a parking spot, and online dating. Optimal stopping problems have attracted considerable research effort, which has primarily focused on showing that optimal policies have intuitively appealing structure.

In optimal stopping problems, the system evolves as a (possibly non-stationary) Markov chain on a set of states S' with transition probabilities $b_n(j|s)$ for $s \in S'$ and $j \in S'$ at epoch n . If the decision maker decides to “stop” in state s at decision epoch n , a reward of $g_n(s)$ is received. If the decision maker decides to “continue,” the decision maker incurs a cost $f_n(s)$. In the finite horizon case, when the problem terminates

after N decision epochs, the decision maker receives a reward $h(s)$ if the Markov chain is in state s at epoch N .

3.8.1 Model formulation

Decision epochs: As noted above, this can be either a finite or infinite horizon model, so

$$T = \{1, \dots, N\}, \quad N \leq \infty.$$

States: The state space is the union of S' and a state Δ that denotes the stopped state:

$$S = S' \cup \{\Delta\}.$$

Actions: Let the action a_0 denote the decision to continue and a_1 represent the stopping decision. Then the action set is

$$A_s = \begin{cases} \{a_0, a_1\} & s \in S' \\ \{a_0\} & s = \Delta. \end{cases}$$

The action to continue in the stopped state is included for completeness.

Rewards: The reward does not explicitly depend on the destination state j , so for $n < N$

$$r_n(s, a) = \begin{cases} -f_n(s) & s \in S', a = a_0 \\ g_n(s) & s \in S', a = a_1 \\ 0 & s = \Delta, a = a_0, \end{cases}$$

with $r_N(s) = h(s)$ for $s \in S'$ if the horizon is finite. In the infinite horizon case, it is more appropriate that $f_n(s) = f(s)$ and $g_n(s) = g(s)$ for all $s \in S'$ and $n = 1, 2, \dots$

Transition probabilities: For $n \leq N$

$$p_n(j|s, a) = \begin{cases} b_n(j|s) & s \in S', a = a_0, j \in S' \\ 1 & s \in S', a = a_1, j = \Delta \text{ or } s = j = \Delta, a = a_0 \\ 0 & \text{otherwise.} \end{cases}$$

Again, in the infinite horizon case, $b_n(j|s) = b(j|s)$ for $n = 1, 2, \dots$

3.8.2 Examples

Selling an asset

A homeowner who is moving to another city has N days to sell a house. Offers arrive throughout the day and by the end of the day, the homeowner has to decide whether to accept the best offer received that day, or wait until the next day for new offers. The set S' represents the set of possible values of the best daily offer for the house, assumed to be around its market value (i.e., bounded) and rounded to the nearest dollar (finite). By waiting until the next day, the homeowner incurs costs $f_n(s)$ related to continued home ownership such as maintenance, mortgage interest, property taxes, and advertising. By accepting the best offer on a given day, the homeowner receives a reward of s , minus the costs associated with selling the house, $L_n(s)$, which includes realtor fees and taxes. Hence $g_n(s) = s - L_n(s)$. At day N , the homeowner must accept the best offer that day, receiving a terminal reward of $h(s) = s - L_N(s)$. The best offer j at decision epoch $n + 1$ is determined by a probability distribution $b_n(j|s)$ that may be conditional on the best offer s in epoch n . This distribution may be non-stationary. For example, early in the planning horizon, rejections could signal that the homeowner expects higher offers in the future than the current best offer. Later in the planning horizon, if bidders know the homeowner must sell, the offers might decrease in value.

Finding a parking spot

A driver seeks a parking spot as close as possible to a restaurant. Assume the driver can move in one direction only and see only one spot ahead. If the spot is not occupied, the driver may choose to park in it or proceed forward. The probability that any spot is unoccupied is assumed to be equal to p , independent of all other spots.

The most natural formulation is as a stationary infinite horizon model with an infinite state space. Elements of S' consist of vectors, (j, k) , where j indicates the location of the parking spot and k indicates whether the spot is free (1) or occupied (0). Let $\mathbb{Z} = \{\dots, 2, 1, 0, -1, -2, \dots\}$ represent the distance of each parking spot to the restaurant where 0 denotes the location of the restaurant, positive numbers denote locations before the restaurant and negative numbers denote locations past the restaurant. Thus, $S' = \mathbb{Z} \times \{0, 1\}$.

States of the form $(j, 1)$ are the only ones in which the action a_1 corresponding to stopping is available. No rewards are received if the individual does not park, that is $f((j, a_0)) = 0$. When the individual parks, the reward is the distance from the parking location to the restaurant, so that $g((j, 1)) = -|j|$.

Transition probabilities of the underlying Markov chain are given by

$$b((j', k')|(j, k)) = \begin{cases} p & j' = j - 1, k = 1 \\ 1 - p & j' = j - 1, k = 0. \end{cases} \quad (3.17)$$

This can be reduced to a finite state formulation by assuming:

- the driver will not start looking for a spot until reaching a distance M before the restaurant, and
- that if the driver has not parked before reaching the restaurant the driver will choose the next available spot.

As a result of the second assumption, the last decision is made in state $(1, 1)$ because the driver will definitely park in state $(0, 1)$ and continue in state $(0, 0)$. Hence the terminal reward is

$$h((0, k)) = \begin{cases} 0 & k = 1 \\ -\frac{1}{p} & k = 0. \end{cases}$$

To see this, note that from 0 onward, the distance to the next available parking spot follows a geometric distribution with parameter p so that

$$|h((0, 0))| = 1p + 2(1-p)p + 3(1-p)^2p + \dots = \frac{1}{p}.$$

Online dating

This example provides a modern take on a classical problem that is often referred to as the *secretary problem*³².

An individual is searching for people to date on a dating app. The dating app shares a brief profile for each potential match, including a photo and description, one at a time. For each potential match, the individual can either “swipe left” to pass or “swipe right” to indicate interest. If the individual swipes left, that profile will not be shown again. If the individual swipes right, a match is made and the two will go on a date. The app offers a free trial until the first match is made or until N profiles have been viewed, whichever comes first. The individual is interested in maximizing the probability of finding the best match during the free trial.

Assume that the N potential matches have an unobservable ranking from 1 to N , with 1 representing the best match. Through the process of examining the profiles, the individual will be able to rank the candidates seen so far relative to each other in a manner that is consistent with the true ranking. If the best profile is seen, the individual will only know that it is the best seen so far, but will not know whether any future profiles will be better. The order of profiles is completely random, so every permutation of the ranks 1 to N is equally likely.

The following formulation may not be immediately obvious, but is the most succinct way to model the problem. Let $S' = \{0, 1\}$, where the state indicates whether the current profile is the best seen so far (1) or not (0). No rewards or costs are accrued if the individual passes by swiping left ($f_n(s) = 0$). In this formulation, rewards correspond to the probability of choosing the best candidate and are only received

³²A closely related problem was first formulated by Cayley [1875] and described in the historical summary in Chapter 1.

upon stopping. If the individual reaches the last profile, the match is automatically made. The terminal reward is thus $h(1) = 1$ and $h(0) = 0$, since the last profile is either the best profile seen so far or not.

If the individual swipes right on the n -th profile, $n \leq N$, and it is not the best profile seen so far, then the probability that the n -th profile corresponds to the best match is $g_n(0) = 0$. But if it is the best profile seen so far, then $g_n(1) = n/N$. To see why this is true, write out the probability that $g_n(1)$ represents explicitly:

$$\begin{aligned} g_n(1) &:= P[\text{profile } n \text{ is rank 1} \mid \text{profile } n \text{ has the highest rank of first } n \text{ profiles}] \\ &= \frac{P[\text{profile } n \text{ is rank 1 "and" profile } n \text{ has the highest rank of first } n \text{ profiles}]}{P[\text{profile } n \text{ has the highest rank of first } n \text{ profiles}]} \\ &= \frac{P[\text{profile } n \text{ is rank 1}]}{P[\text{profile } n \text{ has the highest rank of first } n \text{ profiles}]} \\ &= \frac{1/N}{1/n} = \frac{n}{N}. \end{aligned} \tag{3.18}$$

The second equality follows from the definition of a conditional probability. The third equality is due to the fact that if profile n is the top ranked profile, it must be the top ranked profile within the first n profiles as well. Finally, the fourth equality is due to the fact that the order of the profiles is completely random. So the top ranked profile is equally likely to be in any of the N positions. Similarly, any of the first n profiles is equally likely to be the one with the highest rank.

To determine the transition probabilities, again appeal to the fact that the order of the profiles is completely random. Thus, regardless of the current state, the probability that profile $n+1$ will be the best among the first $n+1$ is $1/(n+1)$. Thus, the transition probabilities are

$$b_n(j|s) = \begin{cases} \frac{1}{n+1} & j = 1, s \in \{0, 1\} \\ \frac{n}{n+1} & j = 0, s \in \{0, 1\}. \end{cases} \tag{3.19}$$

An elegant solution to this problem is provided in Chapter 4.

3.9 Sports strategy

Analytical methods have recently found widespread use in sport decision making, providing many opportunities for applying Markov decision processes. Some applications involve decisions made throughout a game while others are situational. Situational decisions such as whether to “go for it” on fourth down in North American football or whether to steal a base or sacrifice in baseball concern a decision in a particular state in a model of the whole game. They reduce to one-period problems (Section 2.4) when the value function for all games states is estimated from historical data. Other examples such as those described below concern recurrent decisions throughout a game or some portion of it.

3.9.1 When to pull the goalie in ice hockey

In ice hockey, a team has the option of replacing its goalie with an offensive player at any time during a game. This is called “pulling the goalie”. Doing so can be beneficial since there is a greater probability of scoring when an extra offensive player is in the game. However, pulling the goalie also results in a greater likelihood that the opponent scores, since the goal is undefended. Such a goal is referred to as an “empty net goal”. Such a strategy is often employed late in a game when the team is behind, in order to achieve a tie and send the game to overtime.

This decision problem is naturally modeled in continuous time but in keeping with the development of the book, a discrete time formulation follows. Assume that every h seconds a decision is made whether or not to pull the goalie, and that such a decision is not considered prior to M seconds remaining in the game. Usually, M is on the order of 180 seconds (3 minutes).

For concreteness, assume that Team A trails Team B by g goals. Let p_A (p_B) denote the probability Team A (Team B) scores one goal in an interval of length h when Team A pulls its goalie and let w_A (w_B) denote the probability Team A (Team B) scores a goal in an interval of length h when Team A does not pull its goalie. Naturally, $p_A > w_A$ and $p_B > w_B$. It may also be the case that no team scores during an interval of length h , so $p_A + p_B < 1$ and $w_A + w_B < 1$. As is customary when discretizing continuous time models, assume that h is sufficiently small so the likelihood of scoring more than one goal in that interval is negligible.

Decision epochs: Because decisions are made every h seconds up to M seconds before the end of the game,

$$T = \{1, 2, \dots, N\},$$

where $N = M/h$. The first decision epoch corresponds to M seconds left in the game, the second corresponds to $M - h$ seconds left in the game, and so on.

States: The state represents the goal differential, defined as Team B’s goals minus Team A’s goals. Let G be the maximum goal differential at which the coach would consider pulling the goalie. Then

$$S = \{0, 1, \dots, G + 1\}.$$

This formulation adds the absorbing state $G + 1$ to ensure a finite state space. If the goal differential reaches $G + 1$ the coach of Team A will not consider pulling the goalie anymore. If the score differential returns to G because Team A scored a goal without its goalie pulled, the decision problem starts anew. In practice, $G = 3$. A team would not pull its goalie if it is leading, so negative values are omitted from the state space. The state 0 corresponds to a tie score, which also is an absorbing state and the objective of pulling the goalie.

Actions: In all states the coach has the option to not pull the goalie (action a_0) or to pull the goalie (action a_1) when the goal differential is between 1 and G . Thus

$$A_s = \begin{cases} \{a_0, a_1\} & s = 1, \dots, G \\ \{a_0\} & s = 0 \text{ or } G + 1. \end{cases}$$

Rewards: Since the objective is to tie or win by the end of the planning horizon, rewards are only received at termination. So $r_n(s, a) = 0$ for $n < N$ and all s and a . The terminal reward is given by

$$r_N(s) = \begin{cases} 0 & s > 0 \\ 1 & s = 0. \end{cases}$$

Transition probabilities: Under the assumption that one event can occur in a time interval of length h , the transitions probabilities satisfy

$$p_n(j|s, a) = \begin{cases} w_A & j = s - 1, s = 1, \dots, G, a = a_0, \\ w_B & j = s + 1, s = 1, \dots, G, a = a_0, \\ 1 - w_A - w_B & j = s, s = 1, \dots, G, a = a_0 \\ p_A & j = s - 1, s = 1, \dots, G, a = a_1, \\ p_B & j = s + 1, s = 1, \dots, G, a = a_1, \\ 1 - p_A - p_B & j = s, s = 1, \dots, G, a = a_1, \\ 1 & j = s = 0, G + 1, a = a_0, \\ 0 & \text{otherwise.} \end{cases}$$

Application challenges

The key parameters in this model are the relative scoring probabilities. They may vary by team and also depend on the whether the opposing team has been assessed a penalty so that pulling the goalie results in a “two-man advantage” and an increased scoring probability.

The following data comes from the (North American) National Hockey League for the 2013-2020 seasons. When both teams are at full strength (no player is in the penalty box), teams score goals at the rate of 2.25 goals per 60 minutes. When a team pulls its goalie, its scoring rate increases to 6.39 goals per 60 minutes. However, the opposing team’s scoring rate increases to 19.16 goals per 60 minutes. Assuming $h = 5$ seconds, these goal scoring rates correspond to $w_A = w_B = 0.003125$, $p_A = 0.008875$ when pulling the goalie and $p_B = 0.02661$ for an empty net goal. On average, teams pull their goalie around 4 minutes prior to end of the game with a three goal deficit, 2.3 minutes with a two goal deficit, and 1.4 minutes with a one goal deficit. The success rate for pulling the goalie with a one goal deficit is about 14%.

Penalties play a large role in ice hockey. Pulling the goalie when Team B is penalized³³ significantly affects the scoring probabilities, increasing p_A and decreasing p_B relative to the previous values. The same data shows that pulling the goalie when the opposing team has one penalized player increases Team A's scoring rate to approximately 12 goals per 60 minutes ($p_A = 0.0167$), and decreases Team B's scoring rate to approximately 11 goals per 60 minutes ($p_B = 0.0153$).

3.9.2 A tennis handicap system

This section proposes a handicapping system for tennis. Consider a tennis match between two players of unequal skill level. In order to have a fair (and enjoyable) match, the stronger player (Player B) offers the weaker player (Player A) a handicap. The handicap takes the form of a budget of “credits” that Player A can use to win a point without playing it. To our knowledge, such a system has been yet to be widely applied but conceptually presents an interesting strategic challenge: when should Player A use these credits?

Before describing how such a handicapping system may be employed, a basic understanding of tennis scoring is needed. A tennis match consists of games and sets. A player wins a game by scoring four points first, provided that the player “wins by at least two points.” That is, if both players have scored three points, the winning score needs to be at least 5-3. A player wins a set by being to first to win 6 games, again with a “win by two game” rule in effect. If the set score reaches 6-6, then a tiebreaker is played, with the winner winning the set by a score of 7-6. Finally, a match is typically the best two out of three sets or best three out of five sets.

A key feature of this scoring system is that it is “hierarchical”: the match score decomposes into sets and games, each with its own scoring system. Thus the score with Player A serving may be 1-1 in sets, 5-3 in games and 3-1 (commonly referred to as 40-15) in the current game. Faced with this situation, Player A could use a credit to win the game, and hence the set and the match.

Unlike sports such as soccer, in which every goal contributes equally to the final score, not every point is equally valuable in tennis. In fact, a player could win more than 50% of the total points in the match, but still lose the match, due to the hierarchical nature of the scoring system.

To make this precise, if Player A uses a credit at the start of a point, then Player A wins the point, the handicap budget is decremented by one, and the players proceed to play the next point. If Player A decides not to use a credit, then the point is played as usual. Let p_1 (p_0) be the probability³⁴ that Player A wins the point on serve (return) if it the point is played out. Assume that the budget applies to the entire match and that Player A can use a credit regardless of which player is serving.

³³The consequence of a penalty is that one or more fewer players are “on the ice”.

³⁴In most matches the probability of winning a point on serve is greater than when the opponent is serving.

Decision epochs: The start of each point is a decision epoch. Given the scoring system described above, the horizon is infinite but with a random stopping time corresponding to the end of the match. Thus

$$T = \{1, 2, \dots\}.$$

States: The state comprises the current match score, q , the budget of credits remaining, b , and an indicator for the serving player, k . Suppose the set of possible match scores is $Q = \{q_1, q_2, \dots\}$, the starting number of credits is B , and $k = 1$ (0) indicates that Player A (Player B) is serving. The two absorbing states, W and L correspond to Player A winning and losing the match, respectively. Thus, the state space is

$$S = (Q \times \{0, 1, \dots, B\} \times \{0, 1\}) \cup \{W, L\}. \quad (3.20)$$

Each state q represents a set score, a game score and a within-game score. Note that B is at most 24 times the number of sets³⁵ to win the match. In a best two-out-of-three sets match, $B \leq 48$, and in a best three-out-of-five sets match $B \leq 72$.

Actions: Actions are to use a credit (a_1) or not (a_0) at each decision epoch when there are credits remaining. Once no credits remain, the only action is a_0 .

$$A_s = \begin{cases} \{a_0, a_1\}, & s = (q, b, k) \text{ with } b > 0 \\ \{a_0\}, & \text{otherwise.} \end{cases} \quad (3.21)$$

Rewards: Since Player A's objective is to allocate handicap credits so as to maximize the probability of winning the match, the only non-zero reward is when there is a transition to state W , in which case the reward equals 1. Thus

$$r(s, a, s') = \begin{cases} 1, & \text{if } s \neq W, a \in A_s, s' = W \\ 0, & \text{otherwise.} \end{cases} \quad (3.22)$$

Similar to the online dating application, to maximize the probability of an event, the model formulation should be set up so that decision maker receives a reward of 1 only when that event occurs. This follows directly from the observation that the expected value of an indicator variable equals the probability of the indicated event.

Transition probabilities: If a credit is used, then there is a deterministic transition to the state in which Player A has one extra point. Using such a credit could also affect the game and set score, and even result in winning the match. Otherwise, the transitions follow the point-winning distribution of Player A against Player B. For convenience, let q^+ (q^-) denote the score if Player A wins (loses) the point when the

³⁵Since it requires four points to win a game, using 24 credits would allow a player to win a set.

current score is q . Similarly, let k_q^+ (k_q^-) denote the server if Player A wins (loses) the point when the current score is q and the current server is indicated by k . The server changes only if by using the credit, Player A wins the current game. Thus, given that the current state is $s = (q, b, k)$,

$$p(j|s, a) = \begin{cases} 1, & \text{if } j = (q^+, b - 1, k_q^+), a = a_1 \\ p_k, & \text{if } j = (q^+, b, k_q^+), a = a_0 \\ 1 - p_k, & \text{if } j = (q^-, b, k_q^-), a = a_0 \\ 0, & \text{otherwise.} \end{cases} \quad (3.23)$$

Application challenges

The primary challenge in the application of this model revolves around the estimation of the point-win probabilities p_0 and p_1 . Such probabilities depend on several factors including the strength of the opponent (perhaps considering the specific opponent and previous head-to-head successes), the court surface, the weather, and recent playing history. These probabilities are likely to be non-stationary as well due to factors like fatigue or injury.

3.10 The art of modeling

One learns to formulate Markov decision processes by studying how others have done so and trying it out for themselves. Formulations that appear particularly crisp are likely the result of numerous iterations of formulating and re-formulating the problem. By being exposed to and working through many different examples, one starts to build an intuition for how certain problem types are formulated. Below, a systematic approach to model formulation is presented.

How to formulate a Markov decision process model

- **Clearly define the problem.** Verbally describe what the decision maker wishes to achieve, what information is available on which to base decisions, how the system responds to these decisions, and what rewards or costs are incurred as a consequence of the decisions taken. A precise problem description facilitates identifying all model components. Often, however, one must return to, revise or redefine problem characteristics to ensure that the Markov decision process model properly represents the specified situation. The examples above present problem statements that are self-contained, with the information needed to fully formulate the problem.
- **Draw a timeline of events.** Carefully specify when the **state** information becomes available, when **actions** are chosen (i.e., the **decision epochs**), when

rewards are received and when **transitions** occur. Changes to the timing of events can impact the specification of certain model components. Several examples in this chapter illustrate the sequence of events in a typical period. Selected problems at the end of the chapter ask you to reformulate models under modified assumptions about the timing of events.

- **Identify decision epochs.** Specify the precise time at which actions are chosen, using the timeline as a guide.
- **Determine the planning horizon.** Applications may have a horizon that is finite with fixed length, finite with variable length or infinite. Variable length models arise when the policy or realization of a probability take the system to an absorbing state such as in the lion hunting model (Section 3.5), liver transplant model (Section 3.6), the Gridworld model (Section 3.2) and the optimal stopping models (Section 3.8). Most infinite horizon models may be transformed into finite horizon problems through an appropriate reformulation or simply through truncation. An example is provided in the parking problem (Section 3.8.2).
- **Identify states.** The main challenge when formulating a Markov decision process is determining an appropriate state space. Doing so requires taking into account all other model components. Hence, this is the most important step. A well-defined state space can make the rest of the formulation appear obvious; a poorly defined state space may result in complicated (non-Markovian) dynamics, extra computational burden, or insufficient information to write down a complete model. The advance appointment scheduling (Section 3.7) and online dating application (Section 3.8.2) illustrate two challenging examples of state space formulation.

States should encapsulate all of the information available to the decision maker (and no more than is necessary) to specify actions, rewards and transition probabilities. Do not include actions as part of the state space except when one needs to know the current action to decide on future actions³⁶. Often, a time component is required for decision making, but the decision epoch itself may be sufficient to avoid including an extra variable in the state space. Some models may require the addition of zero-reward absorbing states to account for early or random termination times.

- **Specify actions.** Be sure to note whether different sets of actions may be available in each state. Since the Markov decision process formulation requires specifying an action for each state, in *non-actionable* states, that is where there is no meaningful action choice such as an absorbing state, the set of actions should be specified as a single element representing the “do nothing” action.

³⁶As an example consider the queuing service rate control model when there is a fixed cost for changing an action. In this case one needs to know the current action to determine if a switching cost applies.

- **Determine rewards.** Rewards may be stationary or vary with decision epoch. In finite horizon models, be sure to specify a terminal reward function. Also, note whether rewards depend on the subsequent state. It may be possible to model a problem both ways, with rewards that do or do not depend on the next state. However, usually one of these two is more natural. Some of the above examples use a reward function that does not depend on subsequent state. When using criteria based on expected reward, rewards that depend on subsequent states may be replaced by their expectations.

In cases when the decision maker seeks to maximize the probability of an outcome, such as surviving or winning, specify a reward of zero in all states not corresponding to that outcome and a reward of one when that outcome occurs. The reason for this is that the expected value of an indicator of an event equals the probability of the event occurring.

In some reinforcement learning formulations there is no explicit reward function available or its realization is delayed far into the future. In episodic models, the reward may correspond to winning a game or reaching a target state so that the reward is only learned when the episode is completed. In such cases it may prove useful to modify the reward function to indicate progress towards a goal.

While a Markov decision process is generally concerned with maximizing rewards, its formulation, and the reward function specifically, should be independent of the specific choice of optimality criterion such as expected total reward, expected discounted reward, long-run average reward, or expected utility. Also, note that in many applications a decision maker may seek to minimize costs, which can be regarded as negative rewards. Since the Markov decision process formulation here seeks to maximize rewards, costs are best regarded as negative rewards.

- **Specify transition probabilities.** These are often quite complicated and contain many special cases. It is important to appeal to the timeline and the order of events when writing down the transitions. Challenges include taking into account “edge cases” at state space boundaries and noting that some components of the state may evolve deterministically, while others may evolve stochastically. In the presence of absorbing states, be sure to note that under the “do nothing” action the system remains in that state with probability one. For completeness, be sure to note zero probability transitions corresponding to impossible combinations of states and actions, which can be captured under the catch-all heading “otherwise”.

Recall that a Markov decision process with a fixed policy results in a Markov reward process that evolves over a Markov chain. Drawing a Markov chain with directed arcs indicating transitions and denoting probabilities on arcs is a simple but effective method to help ensure that all transitions are accounted for (e.g., probabilities leaving a state for a given action sum to 1) and that they make sense (e.g., transitions occur between states as described in the problem statement).

With complicated multi-dimensional state spaces, drawing such a picture is often a must. Figure 3.7 provides an example.

- **Estimate model parameters.** Most of the examples in this chapter are abstracted from real problem situations. To apply the models in concrete settings, one must estimate the model parameters such as transition probabilities and rewards.

In some cases, such as inventory control (Section 3.1), revenue management (Section 3.3) and queuing control (Section 3.4) the transition probabilities may be derived from parametric distributions with the parameters estimated from historical data. When there are no parametric forms for transition probabilities, care must be taken in estimating probabilities because some may be non-zero but very small. In applications such as the lion hunting model (Section 3.5), clinical decision making (Section 3.6) and sports strategy (Section 3.9) one may appeal to the data and literature from those fields to obtain parameter estimates.

Another challenge when applying models in practice is specifying rewards. In applications where rewards refer to concrete monetary values, specifying rewards can be relatively straightforward. In examples such as lion hunting (Section 3.5), optimal parking and online dating (Section 3.8.2), pulling the goalie (Section 3.9.1), and tennis handicapping (Section 3.9.2) the reward is implicit in the chosen model objective. In other cases, rewards may be derived in consultation with the decision maker.

Bibliographic remarks

Inventory models (Section 3.1) date back at least to Arrow et al. [1951] and Dvoretzky et al. [1952]. The book of Arrow et al. [1958] is an important early reference. Porteus [2002] provides an overview of the historical development in his book on inventory models. Much current research in inventory theory is subsumed under the heading “supply chain management”. Section 8.9.2 of Puterman [1994] provides an inventory example of a constrained MDP with a service level constraint.

The newsvendor model seems to originate with Edgeworth [1888] where it is developed to determine optimal cash reserves to meet random withdrawals from a bank. Arrow et al. [1951] derived the critical fractile solution. It is now described in all operations management textbooks; Chen et al. [2016] provide a modern survey.

The model in Section 3.3 is in the spirit of Gallego and van Ryzin [1994] who provide one of the first examples of dynamic pricing. The book of Talluri and van Ryzin [2004] provides a comprehensive overview of revenue management. Dynamic pricing combined with overbooking has been applied extensively in the airline industry, where it is referred to as yield management Smith et al. [1992].

The queuing control models in Section 3.4 have mostly been studied in continuous time. Early references include Yadin and Naor [1967], Heyman [1968], Naor [1969]

and Sobel [1969]. Our formulation of the discrete time service rate control model follows de Farias and Roy [2003] where they use it to illustrate approximate dynamic programming methods.

The lion hunting example in Section 3.5 is adopted from Clark [1987]. That paper provides many of the parameter values described in the section on application challenges, including the energy storage capacity, daily energy depletion, biomass yield and catch probabilities of gazelles and zebra. The estimate of energy expenditure while hunting was based on Hubel et al. [2016]. Edible biomass of other prey listed were taken from Smuts [1979]. Other applications in ecology include Mangel and Clark [1986], Kelly and Kennedy [1993] and Sirot and Bernstein [1996].

Numerous applications of using Markov decision processes in clinical decision making have appeared in the literature. The model described herein is based on Alagoz et al. [2007], who focus on liver transplantation. The discussion around application challenges is based on methods they employed to specify their model and estimate parameters. Other examples of clinical decision making using Markov decision processes include Shechter et al. [2008], who consider HIV therapy, and Kurt et al. [2011], who model statin treatments for diabetes patients.

The formulation of the advance scheduling model in Section 3.7 follows Patrick et al. [2008]. The result about an optimal policy being independent of the booking window if the window is longer than the largest wait time target and if the system has access to unlimited appointment diversion is given in that paper. Sauré et al. [2012] and Gocgun and Puterman [2014] analyze variants of this model. An example of using the impact of delayed treatment to quantify the cost of delayed imaging appointments appears in Sauré et al. [2012] in the context of radiation therapy.

A Gridworld model appears in Sutton and Barto [2018]. Such models have been widely used in the computer science community to illustrate Markov decision process and reinforcement learning concepts.

Optimal stopping problems originate with early work of Wald [1947], Wald and Wolfowitz [1948] and Arrow et al. [1949]. Karlin [1962] proposes and solves the asset selling problem. The optimal parking problem appears in Chow et al. [1971]. The online dating (i.e., secretary) problem was first proposed by Cayley [1875] in the context of evaluating a lottery.

Sports applications have appeared broadly. The path-breaking monograph of Howard [1960] introduces many of the key Markov decision process concepts, and contains an example of using Markov decision processes in baseball strategy. Carter and Machol [1971] and Chan et al. [2021] develop value functions in football. The pulling the goalie model in Section 3.9.1 originates with Morrison [1976]. A dynamic programming formulation was provided in Washburn [1991]. Hall [2020] provides the recent data quoted in the section on application challenges. The tennis handicapping model in Section 3.9.2 appears in Chan and Singal [2016]. The amazing book by Kemeny and Snell [1960] includes a Markov Chain model for a tennis game.

Exercises

1. Formulate a periodic inventory control problem in which orders arrive after demand has been fulfilled. Clearly show how the timing of events in Figure 3.1 changes.
2. Formulate a periodic inventory control problem in which sales are lost if demand exceeds supply in a period.
3. Formulate a periodic inventory control model in which there is limited storage capacity, limited backlogging and a bound on order size. Assume if the quantity backlogged exceeds its bound, an extra cost incurred.
4. Formulate the reward function of the inventory management example when the revenue of the inventory sold is included.
5. **A newsvendor model with one replenishment opportunity.** Formulate the following variant of the newsvendor model as a Markov decision process. Assume at the start of the period that the newsvendor purchases units of a product from the supplier at a cost of c_1 per item but can purchase additional units at a pre-specified later time during the period at a cost c_2 with $c_2 > c_1$. Such a problem arises when sales exceed initial expectation early in the period so that it may be opportune for the newsvendor to purchase additional items halfway through the period.
6. Formulate a service rate control queuing model with a fixed cost C for changing the service rate. Note that as alluded to above in Section 3.10 it is not sufficient to just modify the reward function.
7. Formulate a combined admission and service rate queuing control problem as a Markov decision process.
8. **Call center routing.** Consider a call center with monolingual (English-only) and bilingual (English and French) call takers. Assume there are two call takers of each type. English-speaking and French-speaking customers call in to the center and indicate their preferred language. English-speaking customers can be served by either type of call taker, but French-speaking customers can only be served by a bilingual call taker. Every minute an English-speaking customer calls in with probability p_E and a French-speaking customer calls in with probability p_F , where $p_E + p_F < 1$. Assume the probability of more than one caller per epoch is negligible. Customers incur a cost of C for every minute spent waiting before service. The duration of a call is geometric with parameter q , regardless of the language. Formulate this routing control problem as a Markov decision process.
9. **Control of a tandem queuing network.** Consider a discrete-time queuing system composed of two single-server queues in tandem and two types of jobs.

Type 1 jobs arrive at queue 1 and after completing service by server 1 also require service by server 2. Type 2 jobs arrive directly to queue 2 and require service at server 2 only. In each period, no job arrives or either a type 1 job arrives, a type

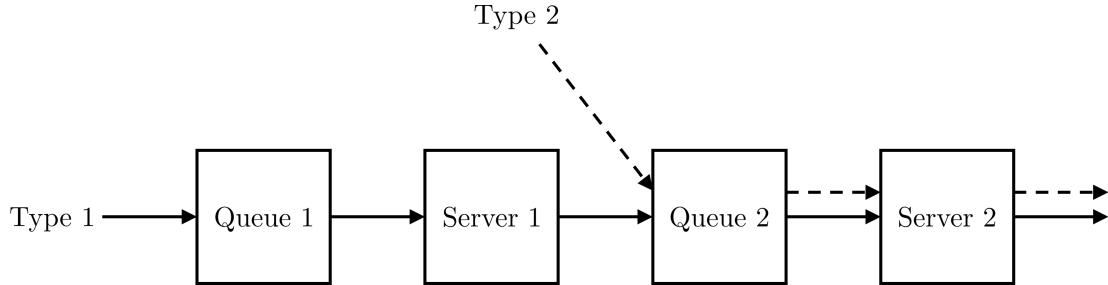


Figure 3.10: Tandem queuing network

2 job arrives or one of each type arrives. Assume the probability of an arrival of a type i job is p_i independent of the whether or not the other type job arrives.

Assume a finite buffer (waiting room) of size M_i in front of queue i . When the buffer is full jobs are blocked (lost) at penalty cost c_i . Completion of a type i job yields revenue R_i with $R_1 > R_2$. In addition, assume a holding cost of $h(s_1, s_2)$ when there are s_i type i jobs in the system (either in the queue or in service).

Formulate the following revenue maximization problems as Markov decision process.

- (a) Job selection: In each period the controller of queue 2 can choose whether to serve a type 1 or type 2 job. Assume that the service is completed with probability q_2 independent of job type and whether or not a job at queue 1 has completed service with probability q_1 . To simplify the formulation assume that if the service is not completed in the current period, the job reverts to the queue prior to the start of the subsequent period.
 - (b) Service rate control: In each period the system can choose the service probability at queue 1 from the set $\{q_{1,1}, q_{1,2}, \dots, q_{1,N_1}\}$ with cost $f_1(q)$ that is non-decreasing in q .
 - (c) Describe and formulate other possible control problems that can apply to this configuration.
10. Formulate a version of the Gridworld navigation model in which there is a positive probability that the robot drops the coffee cup, which is larger if the cup is full. Assume the cup is breakable and if it is dropped, the robot needs to return to the office to retrieve another one.
- Clearly state any assumptions you are making in formulating this model.
11. Formulate a finite horizon version of the Gridworld problem in which if the robot does not return with coffee after N decision epochs, the mathematician gets his

or her own coffee and incurs a penalty of C units. How should C be related to X and R ?

12. **Equipment maintenance.** Formulate the following maintenance problem. You own a piece of equipment that deteriorates over time. While it is operating, it contributes revenue of r dollars per month. When it has been operating for i months since its last maintenance, the probability it fails in the current month is $p(i)$ and the probability it does not fail is $1 - p(i)$. If it fails at any time during a month the cost of repairing it is c_A and it is available for use at the start of the subsequent month. Assume that if it fails during a month, no revenue is generated during that month. On the other hand, the maintenance manager can schedule preventive maintenance in a month at cost c_B . Assume preventive maintenance is always scheduled at the beginning of the month, starts in the first week of the month and takes one month.
 - (a) Draw a timeline for the decision problem and clearly state any assumptions you are making.
 - (b) Formulate the maintenance manager's problem as a Markov decision process. Be clear to state any assumptions you make.
 - (c) In a real application, how do you think $p(i)$ will vary with i and what would be the relationship between c_A and c_B ?
 - (d) Propose a "real life" application of this model.
13. **Stengos and Thomas [1980].** Consider the following generalization of the previous problem. You own two machines which sometimes require maintenance that takes three weeks. Maintenance on one machine costs c_1 per week while maintenance on two machines costs c_2 per week. The probability either piece of equipment breaks down if it has operating for i weeks is $p(i)$. Assume that if the equipment breaks down during a week, maintenance begins at the start of the next week. However, if you decide to perform preventive maintenance, you do so at the start of a week. Moreover, assume that the two pieces of equipment break down independently.
 - (a) Formulate this problem as an infinite horizon Markov decision process.
 - (b) How are c_1 and c_2 related? Why?
14. Reformulate Exercise 13 assuming that the time it takes to complete maintenance on a piece of equipment is random. In particular, assume that maintenance is completed in any period with probability q , independent of other periods.
15. **Bertsimas and Shioda [2003].** A restaurant contains both two-seat and four-seat tables. Parties of two and four arrive randomly and request service. No reservations are taken. When a party of four arrives and a four-seat table is

empty, they should be seated but should the manager ever seat a party of two at a four seat table? If so, when? Also, when should requests for service be denied, if ever?

Formulate this problem as Markov decision process assuming the following, unrealistic as it may be. Decisions are made every 10 minutes and the restaurant operates 24 hours a day. There are two two-seat and two four-seat tables. Meals consist of two courses, durations of each are geometrically distributed independent of party size. Course 1's completion probability per epoch is 0.7, while course 2's is 0.8. In any period there is at most one arriving party. A party of two arrives with probability 0.2 and a party of four with probability 0.1. Assume that the waiting area holds at most 6 people. If it is full, arrivals are blocked and do not enter. Also any waiting party may leave in a 10 minute period with probability 0.05.

Revenue is as follows. A party of 2 contributes \$50 and a party of 4 contributes \$100. The cost of waiting (incurred by the restaurant) is \$6 per person per hour.

16. Formulate the advance scheduling problem when appointments not booked on first day available are added to the next days demand with cost C' .
17. Formulate the advance scheduling problem where instead of the target representing a fixed day, target windows are used for each urgency class. Let T_k^l and T_k^u be the lower and upper limits of the target window for urgency class k . If an appointment is scheduled within this window, no costs are incurred. If a class k appointment is scheduled before (after) T_k^l (T_k^u), then a cost of C_k^l (C_k^u) is incurred.
18. Modify the lion hunting problem to take into account that on any day, the lion may be captured by poachers with probability $1 - \lambda$. How is this related to discounting?
19. At the start of each day, a lion decides whether or not to hunt and if so, in what group size. The probability of catching prey varies with group size. This presents a trade-off: a larger group has a greater probability of a successful hunt, but then less food is available for each lion in the group. Assume a maximum group size of M and that all captured prey is split evenly among the group. Let λ_m , $m = 1, 2, \dots, M$, denote the probability that a group of size m is successful in its hunt. Assume the prey being hunted yields a total edible biomass of e units. Thus if the lion hunts in a group of size m and is successful, it receives e/m units of edible biomass.
20. Reformulate the original lion hunting behavior model so that the objective is to maximize the number of days of survival instead of the probability of survival. Clearly note what changes are necessary.

21. **The ride-sharing driver's dilemma.** At random times throughout the day, a ride-sharing driver receives offers of potential trips, including their expected revenue and time to complete the trip. The driver can either accept the trip or decline it and wait for the next offer. Formulate this problem as a discrete time Markov decision process clearly stating all assumptions being made.
22. Consider a variant of the online dating problem in which the decision maker's goal is to maximize the probability of choosing one of the two best candidates. How would you modify the formulation to take this into account.
23. Consider a variant of the online dating problem in which the decision maker's goal is to maximize the rank of the selected date. Modify the formulation accordingly.
24. Reformulate the tennis handicapping problem to account for second serves. That is, if a player misses a first serve, they have a second chance to get the ball in before losing the point. Suppose the probability that the first serve goes in is q_1 and the probability the second serve goes in is q_2 . Conditioned on the serve going in, the probability of winning the point is $p_{1,1}$ and $p_{1,2}$ for the first and second serve, respectively. Since first serves tend to be more aggressive, $q_1 \leq q_2$ and $p_{1,1} \geq p_{1,2}$. Assume handicap credits can only be used when serving.
25. *"Scrabble, like life, is a trade-off between today and tomorrow-between spending and saving. It's what an economist would call a dynamic programming problem."* [Roeder, 2022] The game of Scrabble provides an opportunity for applying Markov decision processes. Develop a model for a decision of whether a player should replace some or all tiles during a turn. This is a rather complex model that will be challenging to formulate in its entirety. We do not believe it has been addressed in the Markov decision process literature. It may be amenable to reinforcement learning methods.