

Introduction to Markov Decision Processes

Martin L. Puterman and Timothy C. Y. Chan

January 10, 2023

Chapter 1

Examples and Applications

This material will be published by Cambridge University Press as Introduction to Markov Decision Processes 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, 2021.

*We welcome all feedback and suggestions at:
martin.puterman@sauder.ubc.ca and tcychan@mie.utoronto.ca*

*For the things we have to learn before we can do them, we learn by doing them.
Aristotle, Philosopher, 384-322 BC*

Representing a dynamic decision problem as a Markov decision process requires specifying all of the model components described previously: decision epochs, states, actions, transition probabilities and rewards. In this chapter, we show how to do so by identifying these objects in many different examples. We have chosen applications from a broad range of disciplines to illustrate both the wide applicability of the Markov decision process formalism and the many features of model formulation. We encourage the reader to attempt formulating the models before seeing how we did so. We conclude the chapter with a section that provides guidance on how to formulate Markov decision process models. The problems at the end of the chapter provide further opportunities to formulate Markov decision processes or revise the examples presented herein when assumptions differ.

1.1 Revenue Management: Using Price to Manage Demand

This problem describes a different approach to inventory management. It pertains primarily to products that decline in value over time. As a concrete example, consider the challenges faced by our acquaintance, Frank Z., the former owner of a chain of women's fashion stores in Vancouver, Canada. He related to us 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've lost the opportunity to sell my inventory and also incur costs to store it."

Let us formalize Frank's problem. A retailer has an inventory of M units of a product at the beginning of the season and requires a policy to vary prices over the N month season so as to maximize revenue. We assume prices are set at the beginning of each month, to be chosen from a finite set of prices, and cannot be changed during 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 $\lambda_{n,a}$ is non-increasing in n because fashion products may become less trendy as time goes on and expected demand will decrease. It is also reasonable to assume that $\lambda_{n,a}$ is non-increasing in a since in a typical demand curve the quantity demanded decreases as price increases.

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 of 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 price to set in each period. Assuming there are K candidate prices for all $s \in S$,

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

We assume that $a_1 \leq a_2 \leq \dots \leq a_K$ and that $a_1 = H$ denotes the scrap value.

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 $r_N(s) = a_1 s$.

Transition Probabilities: Since no inventory is added during the planning horizon, only transitions to states with ending inventory $j \leq s$ have nonzero probability. If the demand is at most $s - 1$, then the ending inventory is at least 1. If the demand equals or exceeds s , then the ending inventory is 0. This logic leads to the following transition probabilities:

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

Application Challenges

Applying this model presents two challenges: determining the set of mark down prices and estimating the time varying demand function parameters. In retail, the mark down 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 not enough historical data for the product. The parameter $\lambda_{n,a}$ may itself be represented as a function of n and a and learned in the decision problem.

1.2 A Periodic Review Inventory Model

“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 next order should be. Underestimating would hurt not just his sales but the status of his Amazon listing; overestimating would drain cash upfront, and he would incur further charges from Amazon for storing excess inventory in its warehouses.

The Great Amazon Flip-a-Thon; John Herrman, New York Times, S1,S7, April 4, 2021

Inventory models represent some of the earliest and most widely studied Markov decision process models. They concern determining appropriate inventory levels for a retail product in the face of future random arrivals of customer orders. As the quote above alludes to, 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, and this trade-off is the main tension that one needs to balance.

We assume that an *inventory manager* periodically (hourly, daily, weekly or monthly) observes the inventory level of a product and, if deemed opportunistic, places an order

with a *supplier*. The order from a supplier may arrive immediately, by the next review period or after several periods. The delay between placing and receiving an order is referred to as a *lead time*. An inventory model may contain some or all of the following features:

1. **Ordering costs** may consist of a fixed and a variable component. The fixed component K represents the administrative cost of placing an order and the variable component $c(u)$ represents the cost of ordering u units. When no order is placed, the ordering 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 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. Alternatively, demand may be *backlogged*, meaning that it is not lost and will be fulfilled when future inventory arrives.
4. **Penalty costs**, denoted by $p(u)$, represent the cost to the inventory manager when demand exceeds inventory by u units in one period (i.e., due to backlogging). Assume $p(0) = 0$.
5. **Revenue** of $R(u)$ is received when u units are sold.
6. **Customer demand** for units arrives during a period. The demand distribution may be known or unknown, and may be static or time-varying. Let random variable D_n denote the demand in period n . Therefore, the probability that d units are demanded in period n is $P(D_n = d)$.
7. **Product** may have an “infinite” shelf life or may be *perishable*. Perishable items may last one period (for example a newspaper, a loaf of fresh bread, a defrosted vaccine or a seat at a sporting event) or may last for multiple periods (for example blood products, new electronics or fashion goods). Assume product is only available in whole units.
8. **Scrap value**, denoted by $H(s)$, equals the value of the ending inventory when there are s units on hand and the planning horizon is finite. 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).

Formulating this problem as a Markov decision process requires precisely specifying the timing of events. Figure 1.1 depicts the event sequence for the model formulated below. In a typical period, the decision maker reviews the inventory level at a decision epoch and decides how many units of the product to order from the supplier. The

ordered products arrive before the end of the same period. Demand arrives throughout the period and is fulfilled at the end of the period, using inventory on hand as well as the products that arrive from the current period's order. If demand exceeds inventory, it is backlogged for future fulfillment. In the formulation below, we focus on cost minimization, thus ignoring the revenue from selling the inventory, which is left as Exercise 27.

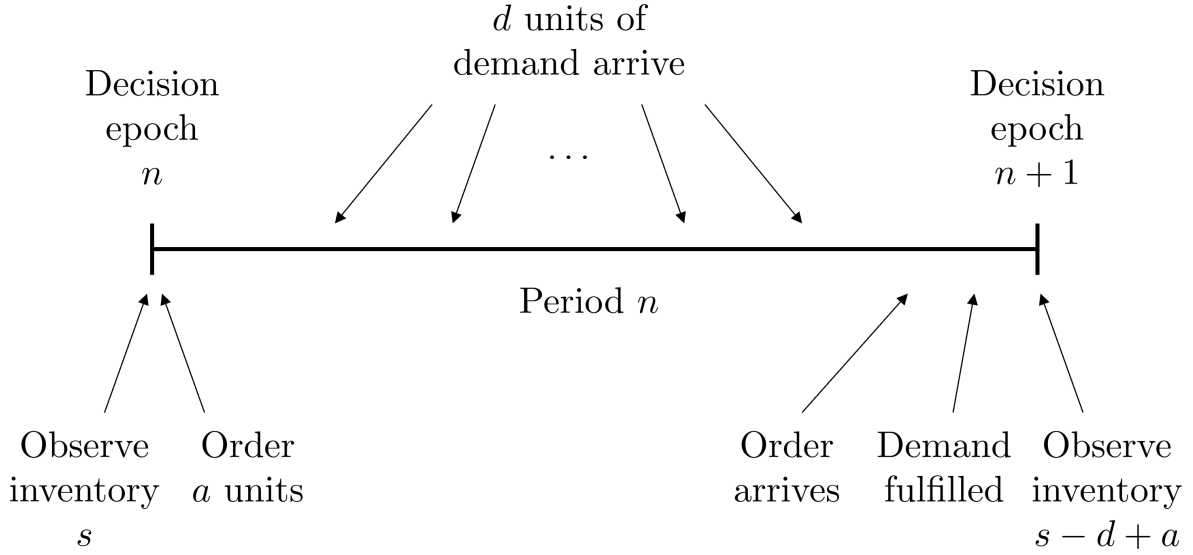


Figure 1.1: Timing of events in the periodic review inventory model described below.

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 ensuing 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 set can be truncated to make it finite.

Rewards: Since our formulation assumes reward maximization, we write the reward function as the negative of the costs. The reward consists of the ordering costs, which are incurred if $a > 0$, and holding or penalty costs if s is positive or negative, respectively. Our simplifying assumption is that holding and penalty costs are assessed at the beginning of the period based on the *starting* inventory.

Let $I(\cdot)$ denote an indicator variable. Then the reward can be written

$$r_n(s, a, j) = -KI(a > 0) - c(a) - h(s)I(s > 0) - p(-s)I(s < 0)$$

when $n < N$. When the state at the next decision epoch $n + 1$ is j , this means that demand in period n was $s + a - j$. The argument of $p(\cdot)$ is $-s$, which equals the backlogged demand when s is negative. If N is finite, $r_N(s) = H(s)$, the scrap value of the remaining inventory at the end of the planning horizon.

Transition Probabilities: Since the state at the next decision epoch cannot exceed $s + a$ (i.e., the demand cannot be negative), the transition probabilities are

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

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 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.

1.3 Discrete-Time Queuing Models

A queuing system consists of arrivals, a queue and one or more servers. Jobs arrive, wait in a queue if the servers are busy, are served by a free server and then depart the system. 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.

We will illustrate the formulation of service rate and admission control in the following subsections. A routing control example is provided in Exercise 6.

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. Here, we consider a discrete-time formulation. 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 . Let the set of decision epochs be denoted by $\{1, 2, \dots\}$ corresponding to times $\{h, 2h, \dots\}$. Assume h is sufficiently small so that it is very unlikely more than one event (an arrival or service completion) occurs during that time interval.
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.

1.3.1 Service Rate Control

Consider a queuing system with a single server and an infinite capacity queue. Each decision epoch, the decision maker chooses a service “rate” $a_k, k = 1, 2, \dots, K$ that denotes the probability a job being processed is completed in the current period. Assume the probability a job arrives between two decision epochs is b , independent of the number of jobs in the system. We assume that $a_1 \leq a_2 \leq \dots \leq a_K$ and $a_K + b \leq 1$. The queuing system is shown in Figure 1.2.

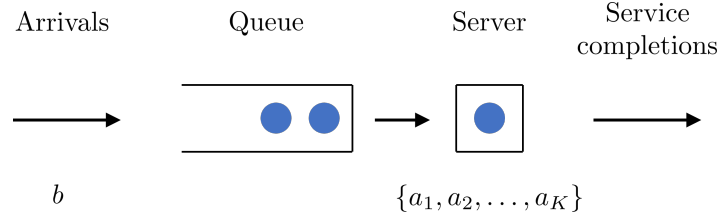


Figure 1.2: Schematic representation of a single server queuing system with adjustable service rate.

Assume 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, such that both $m(a)$ and $f(s)$ are non-decreasing in their arguments.

Decision Epochs:

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

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

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

Actions: Actions represent the probability a job currently being served is completed before the next decision epoch. For $s \in S$,

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

However, given the assumed cost structure, we can assume without loss of generality that $A_0 = \{a_1\}$.

Rewards: Costs may be regarded 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 $k = 1, 2, \dots, K$,

$$p(j|s, a_k) = \begin{cases} a_k & j = s - 1 \\ b & j = s + 1 \\ 1 - a_k - b & j = s. \end{cases} \quad (1.1)$$

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

$$p(j|0, a_k) = \begin{cases} b & j = 1 \\ 1 - b & j = 0 \\ 0 & \text{otherwise} \end{cases} \quad (1.2)$$

The above formulation assumes an unbounded, discrete state space. For computation, we must truncate the state space at a large value, say W , and assume that arrivals are *blocked* when the system state is W . This means that to complete the formulation we must add the additional transition probabilities $p(W - 1|W, a_k) = a_k$ and $p(W|W, a_k) = 1 - a_k$ and limit (1.1) to only $s = 1, 2, \dots, W - 1$.

1.3.2 Admission Control

In an admission control model, the decision maker or “gate keeper” decides whether or not to admit an arriving job into a queuing system with fixed arrival and service rates. This example provides an illustration of a model with non-actionable states, which occur when no job arrives in the preceding period.

Assume at most one job can arrive between decision epochs and it does so with probability b . If it does not get admitted by the decision maker, the job is lost. Let $h(j)$ 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 two decision epochs. When a job is admitted to the system, the system receives a payment of R . In addition, we assume $b + w \leq 1$, which is reasonable when the time discretization step h is small. Figure 1.3 provides a schematic representation of the system and Figure 1.4 depicts the one-period dynamics.

Decision Epochs:

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

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

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

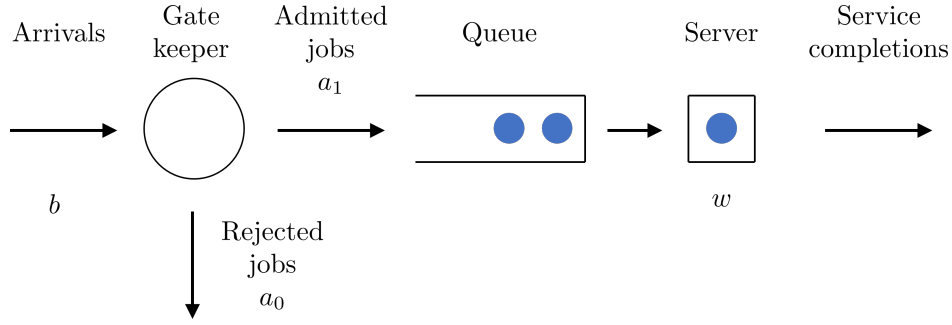


Figure 1.3: Schematic representation of a single server queuing system with admission control.

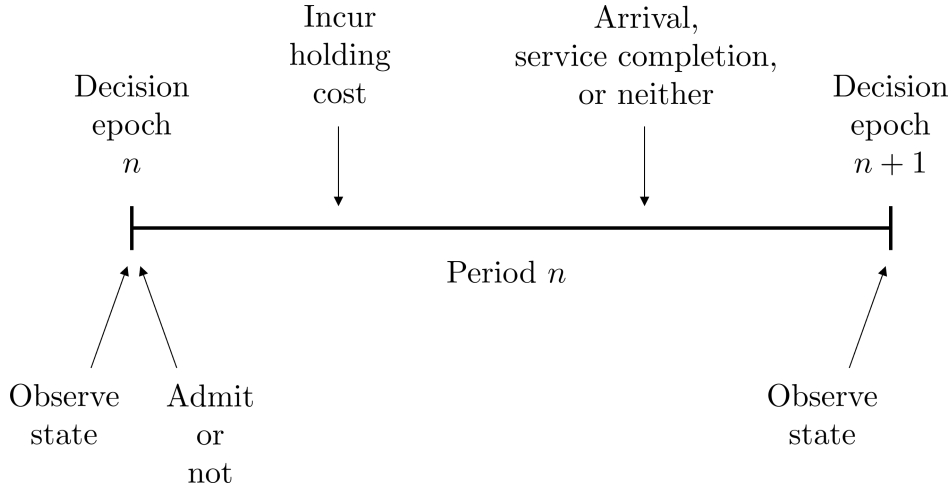


Figure 1.4: Timing of events in the queuing admission control problem.

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

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

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.

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

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

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, for $j = 1, 2, \dots$, $a = a_0$, and $k = 0$ or 1 ,

$$p((j', k')|(j, k), a) = \begin{cases} w & j' = j - 1, k' = 0, a = a_0 \\ b & j' = j, k' = 1, a = a_0 \\ 1 - b - w & j' = j, k' = 0, a = a_0. \end{cases} \quad (1.3)$$

Since service completions are not possible if the system is empty, the above dynamics become

$$p((j', k')|(j, k), a) = \begin{cases} b & j' = j, k' = 1, a = a_0 \\ 1 - b & j' = j, k' = 0, a = a_0, \end{cases} \quad (1.4)$$

when $j = 0$, $a = a_0$ and $k = 0$ or 1 .

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

$$p((j', k')|(j, k), a) = \begin{cases} w & j' = j, k' = 0, a = a_1 \\ b & j' = j + 1, k' = 1, a = a_1 \\ 1 - b - w & j' = j + 1, k' = 0, a = a_1. \end{cases} \quad (1.5)$$

Note the inclusion of $j = 0$ in (1.5). In contrast to (1.3), we can include $j = 0$ into (1.5) 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 1.5 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. We will revisit this example in Chapter ??, provide a simpler formulation of the model in terms of the *post-decision state*, and illustrate some computational results. The post-decision state provides an alternative view of the decision timeline depicted in Figure 1.4 that allows the transition dynamics to be modeled more easily. As discussed at the start of the chapter, drawing a correct timeline is an important part of formulating the model correctly. **(June 21 - check that we do this in chapter 5).**

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

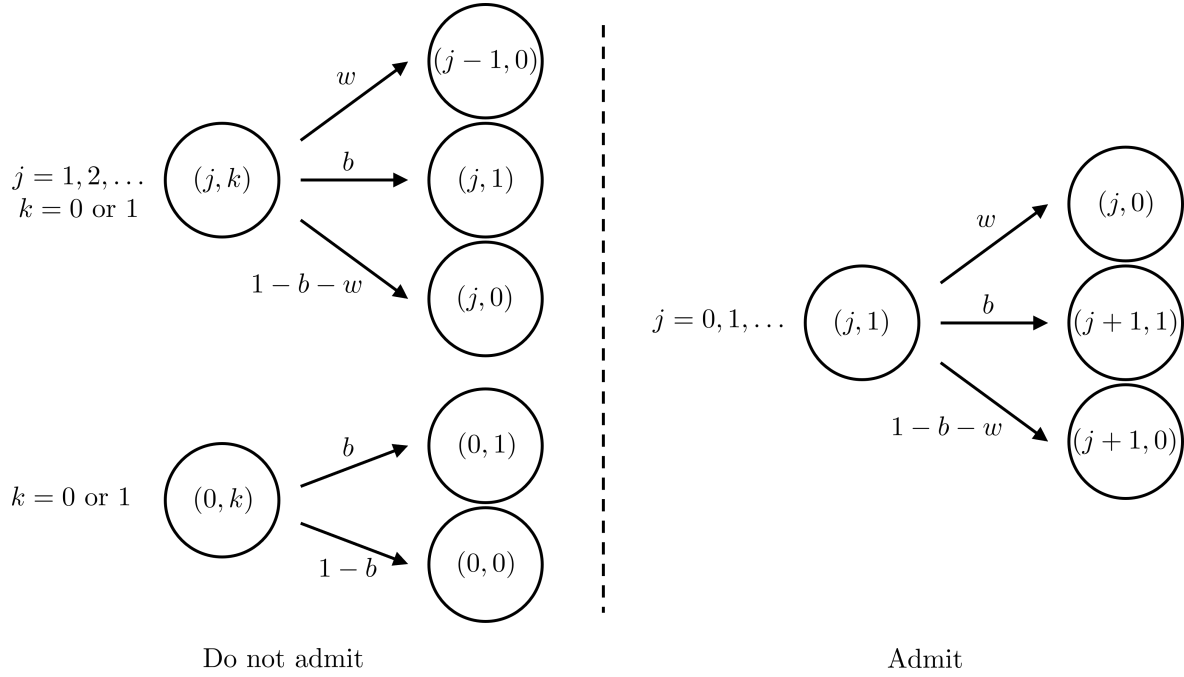


Figure 1.5: Possible state transitions for admission control problem.

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 for 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.

1.4 Lion Hunting Behavior

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 an optimization model can explain observed animal behavior.

As an example, consider the challenge facing a lion (*panthera leo*) when deciding to hunt for food. 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 for impalas, gazelles, wildebeests, giraffes and zebras. About half of the time they hunt in groups. Here we assume that the lion hunts alone. (Exercise 20 asks you to formulate the group size decision problem.) Catch probability may vary 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 ; $m = 1, 2, \dots, M$. We assume a successful hunt for species m yields a total e_m units of energy. For simplicity, we assume all relevant quantities are rounded to the nearest integer.

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. Naturally, $S = [0, C]$, but to maintain a finite state formulation, we discretize the state to the nearest energy unit so that

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

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 ”.

Rewards: Given the lion's survival objective, the lion receives a reward of 1 if it is alive at the end of the season and a reward 0 if it is 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 $\min\{s - h + e_m, C\}$. If it is unsuccessful, 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, we recommend conducting sensitivity analyses of these parameters.

1.5 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 we describe an application to liver transplantation.

Suppose a patient with end stage liver disease requires a liver transplant and a compatible liver becomes available from a recently deceased individual. Depending upon the quality of the offered organ, the patient's medical team may either accept the offered organ and receive a transplant immediately, or reject it and wait for the next offer. For example, a relatively healthy patient may choose to reject a lower-quality liver in the hopes of being offered a higher-quality liver in the future. On the other hand, a patient in poor health may be better off accepting the first liver available. A Markov decision process model can be used to formalize this decision problem and weigh the trade-offs.

To facilitate modeling, patient health status and organ quality are grouped into discrete categories. Patient health states are ordered from 1 to $H+1$ where 1 represents the healthiest state, H represents the least healthy (but still alive) state, and $H+1$ represents death. Similarly, suppose there are L liver quality states, ordered in quality from 1 (highest) to L (lowest). Let state $L+1$ represent the case where no liver is offered.

Rewards measure life expectancy in days. Each day in which no transplant is made, the patient accrues a reward of 1, representing an extra day alive, independent of the health state. The life expectancy post-transplant is $R(h, l)$, if a patient in health state h accepts a liver of quality l . It is reasonable that $R(h, l)$ be non-increasing in h and l .

If a patient does not receive a transplant, the patient's health state either remains the same or deteriorates. Let $b_h(k)$ denote the probability that the health of a patient in health state h deteriorates by $k = 0, 1, \dots, H+1-h$ categories in one period and let $w_h(l)$ denote the probability that a patient in health state h is offered a liver of quality $l = 1, \dots, L+1$ at the start of a period. Assume these distributions are stationary and independent.

Decision Epochs: Given that the unit of reward is a day, we consider daily decision epochs with an infinite planning horizon. A practical upper bound on the number of decision epochs might be 26,000 (assuming no transplants for people over 90 or younger than 20). However, 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.

Thus,

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

States: States represent the patient's health and liver quality, as described above. We also add an absorbing state C to represent the post-transplant state. For convenience, we define $S_H = \{1, \dots, H\}$, $S_L = \{1, \dots, L\}$, and $S_{L+1} = \{1, \dots, L+1\}$. Then,

$$S = (S_H \times S_{L+1}) \cup \{H+1, C\}.$$

Actions: At decision epochs when an organ is available, let a_t represent the action to accept an organ and a_w represent the action to wait, i.e., do not accept an organ. We also let a_w 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.

$$A_s = \begin{cases} \{a_t, a_w\} & s \in S_H \times S_L \\ \{a_w\} & s \in (S_H \times \{L+1\}) \cup \{H+1, C\}. \end{cases}$$

Rewards: The reward function is

$$r(s, a, j) = \begin{cases} R(h, l) & s = (h, l) \in S_H \times S_L, a = a_t, j = C \\ 1 & s \in S_H \times S_{L+1}, a = a_w, j \in \{h, \dots, H\} \times S_{L+1} \\ 0 & s \in S_H \times S_{L+1}, a = a_w, j = H+1 \\ 0 & s \in \{H+1, C\}, a = a_w, j = s. \end{cases}$$

Transition Probabilities: If the patient does not receive a transplant, then a transition from health state h to state h' means that the patient's health deteriorates by $h' - h$ levels. If the patient receives a transplant, then the transition is deterministic to the post-transplant state. Similarly, if the action is to “do nothing” in the post-transplant or death states, then the transition is a deterministic self-transition.

$$p(j|s, a) = \begin{cases} b_h(h' - h)w_h(l') & s = (h, l) \in S_H \times S_{L+1}, a = a_w, j = (h', l') \in \{h, \dots, H\} \times S_{L+1} \\ b_h(H+1 - h) & s = (h, l) \in S_H \times S_{L+1}, a = a_w, j = H+1 \\ 1 & s \in S_H \times S_L, a = a_t, j = C \\ 1 & s \in \{H+1, C\}, a = a_w, j = s \\ 0 & \text{otherwise.} \end{cases}$$

Application Challenges

Application of Markov decision process models to clinical decision making requires 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 and the progression of end stage liver disease. Data required includes suitable 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. For instance, patient health status can be measured using the Model for End Stage Liver Disease (MELD) score, which is a function of various laboratory values. The scores range from 6 to 40, with higher scores indicating poorer health and higher probability of mortality. If data is sparse, MELD scores can be aggregated to form coarser states. A similar approach can be taken when defining liver quality states, which may depend on donor age, race and sex. A range of historical data can be used to calculate other parameters. For example, data from UNOS can be used to estimate $R(h, l)$, via a proportional hazards model, and organ arrival rates. Transitions between health states can be modeled using a natural history model.

1.6 Advance Appointment Scheduling

In many applications, decision makers must allocate scarce resources prior to the arrival of future random demand. As an example, a hospital diagnostic imaging department faces the challenge of scheduling appointments for current medical imaging requests so as to meet clinical wait time targets, without knowing exactly how many and when future requests will arrive.

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 cost 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 the delayed imaging. We assume that $C_1 > C_2 > \dots > C_K$ to represent that the higher delay cost are associated with more 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 formulation, we 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. We 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, a scheduling clerk assigns an appointment date to each request and informs the patient. Figure 1.6 provides a timeline for this process. The challenge is to schedule today's requests before realizing future requests for appointments in the face of limited capacity.

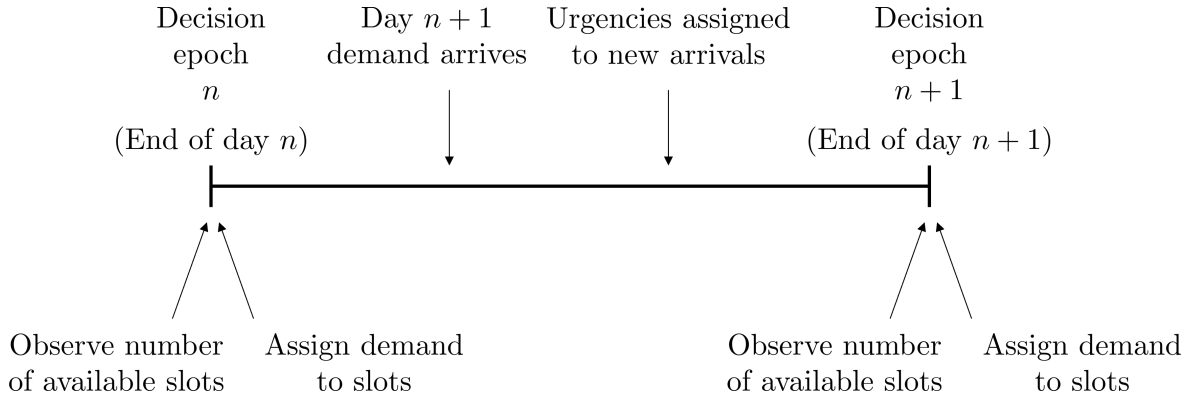


Figure 1.6: Timing of events in the appointment scheduling model.

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?
7. How far in advance can appointments be scheduled?

In our example, we assume: 1) access to overtime, 2) fixed appointment lengths, 3) no no-shows / late cancellations 4) no rescheduling, 5) flexible target dates (but with a penalty for exceeding the target date), 6) stationary arrivals and uncorrelated demand between urgency classes, 7) a fixed appointment *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 15 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 $s = (b_1, b_2, \dots, b_N, w_1, w_2, \dots, w_K)$. Let $b_i \in \mathcal{B} = \{0, 1, \dots, B\}$ denote the number of appointments that have already been booked on day i for $i = 1, \dots, N$. Let $w_k \in \mathcal{W}$ denote the number of appointments of urgency class k waiting to be scheduled at a decision epoch. Note that if there are b_i appointments booked on day i , then there are $B - b_i$ remaining appointment slots on that day.

To simplify notation, we introduce the vectors $\mathbf{b} := (b_1, b_2, \dots, b_N)$ and $\mathbf{w} := (w_1, w_2, \dots, w_K)$. Thus, a specific 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 have already been scheduled. 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 farther out since we assume access to unlimited 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 length NK , K and $(N + 1)K$, respectively. The action set $A_{(\mathbf{b}, \mathbf{w})}$ is

$$A_{(\mathbf{b}, \mathbf{w})} = \left\{ (\mathbf{x}, \mathbf{y}) \geq \mathbf{0} \left| \sum_{n=1}^N x_{kn} + y_k = w_k \text{ for } k = 1, \dots, K \text{ and } b_n + \sum_{k=1}^K x_{kn} \leq B \text{ for } n = 1, \dots, N \right. \right\}. \quad (1.6)$$

The first condition in (1.6) 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: We can write the penalty cost associated with scheduling an urgency class k request n days from the current day as $C_k(n - T_k)^+$. This function specifically models the cost being 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 16 considers the variation where instead of the target representing a fixed day, a target window is used.

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. Notice that the reward does not depend on the next state after (\mathbf{b}, \mathbf{w}) since the updated \mathbf{b} vector is entirely determined by the current actions and the cost of future appointment requests will be captured in the reward of the new state with an updated \mathbf{w} vector.

Transition Probabilities: Transition probabilities depend on both the action choice and the random arrival distribution. At the start of a period, the calendar moves 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:

$$p((\mathbf{b}', \mathbf{w}') \mid (\mathbf{b}, \mathbf{w}), (\mathbf{x}, \mathbf{y})) = \begin{cases} \prod_{k=1}^K p_k(w'_k) & \mathbf{b}' = (b_2 + \sum_{k=1}^K x_{k2}, b_3 + \sum_{k=1}^K x_{k3}, \dots, b_N + \sum_{k=1}^K x_{kN}, 0) \\ 0 & \text{otherwise.} \end{cases} \quad (1.7)$$

Some comments about this formulation follow:

- The concept of a booking horizon may be regarded as an artifact of the modeling process and imposed to maintain a finite state space. We can use a fixed booking horizon since our model considers the availability of unlimited overtime. Without overtime, an unbounded booking horizon may be needed.
- 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.

- 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.
- 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 ??).
- 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 ?? and Chapter ??, it may be more convenient to formulate the model in terms of post-decision states.
- 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 cases are dealt with, 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, for example 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. The above model formulation was based on a real application, and 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.

1.7 Grid World Navigation

“A mathematician is a machine for turning coffee into theorems.”

Alfred Renyi, Mathematician, 1921-1970

A working mathematician frequently requires coffee and, to avoid time away from theorem proving, employs a robot to bring coffee when needed. 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 falling down an open stairwell. Figure 1.7 depicts the grid the robot must navigate to retrieve and deliver coffee.

We assume that the robot knows the arrangement of the grid, which grid cell it occupies and the location of the grid boundaries. This means the robot will never attempt to move outside the grid boundary. Variations of this problem may consider the situation where the robot does not know its location with certainty or the configuration of the grid. Section ?? describes one such example.

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 cardinal 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, it does not move further. We 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 state there are k possible locations where the robot can end up (including the robot staying in the same location), then 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 with probability p_F and it moves left, 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, we assume $p_E > p_F$.

The goal is to return with a full coffee cup at the earliest possible decision epoch. If the robot successfully delivers coffee to the mathematician it receives a bonus reward equivalent to B epochs. If it falls down the stairs it incurs a penalty of X because the mathematician has to interrupt work to rescue the robot. We assume $X \gg B \gg 1$.

Decision Epochs: We 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.

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

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

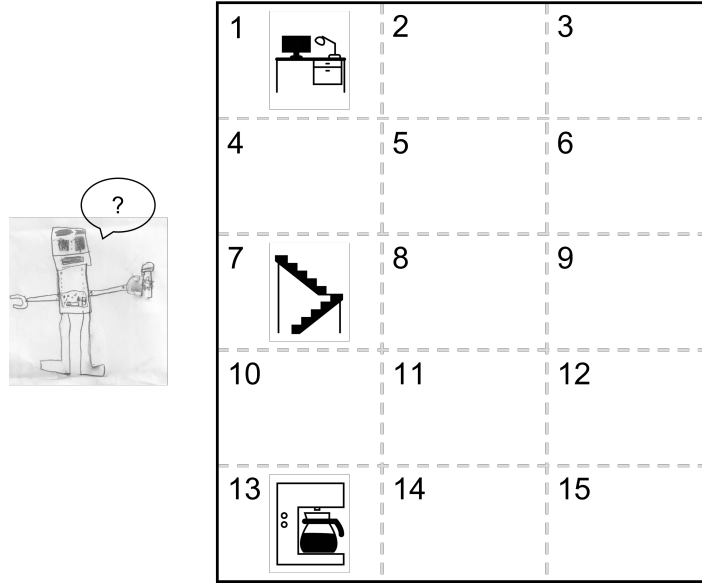


Figure 1.7: Schematic layout for grid world navigation example.

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,

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

For convenience, we will refer to the office as O (grid cell 1), the stairwell as ST (grid cell 7), and the coffee room as CR (grid cell 13).

Actions: Actions represent the direction the robot attempts to travel. Assume the robot can only move north, south, east and west, denoted by U , D , R and L , respectively. Because we assume the robot knows the layout of the grid and its location, the grid boundary constrains its intended movement. For example,

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

We distinguish action sets for some particular states as follows:

$$A_{O,F} = A_{ST,F} = A_{ST,E} = \{a_0\}, \quad A_{CR,E} = \{a_1\}.$$

The states (O, F) , (ST, F) and (ST, 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, we assume it takes one period to fill it. We denote the action of filling the cup by a_1 .

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

$$r((2, F), a, (O, F)) = B - 1, \quad a \in A_{(2, F)} \quad (1.8)$$

Transitions into the stairwell receive a reward of $-X - 1$ regardless of the status of the cup. For example,

$$r((4, k), a, (ST, k)) = -X - 1, \quad a \in A_{(4, k)} \text{ and } k \in \{E, F\}. \quad (1.9)$$

All other feasible state transitions between neighboring cells result in a reward of -1 . For example,

$$r((5, k), a, (6, k)) = -1, \quad a \in A_{(5, k)} \text{ and } k \in \{E, F\}. \quad (1.10)$$

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

$$r((O, F), a_0, (O, F)) = 0 \text{ and } r((ST, k), a_0, (ST, k)) = 0, \quad k \in \{E, F\}. \quad (1.11)$$

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

Transition Probabilities: We provide some typical probabilities:

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

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

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):

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

Some comments about this example follow:

1. This stylized example is in the spirit of numerous problems that have appeared in the reinforcement learning literature on robotic control. It combines features of stochastic shortest path and gambler's ruin problems.
2. 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 modeling uncertainty in intended movement, specifying the behavior of the robot's sensors, and providing the robot with a mapping of the area. These challenges are amplified when the application involves robotic movement in a three-dimensional space.

1.8 Optimal Stopping

An elegant collection of applications belongs to the class of problems known as “optimal stopping problems”. Examples include selling an asset, finding a parking spot, or online dating. We describe these examples after formulating the general model. Optimal stopping problems have attracted considerable research effort, which has primarily focused on showing that an optimal policy has an 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 t , the decision maker receives a reward of $g_n(s)$. 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 .

1.8.1 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 C (for *continue*) denote the decision to not stop and Q (for *quit*) represent the stopping decision. Then the action set is

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

We include the action to continue in the stopped state for completeness.

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

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

with $r_N(s) = h(s)$ for $s \in S$ if the horizon is finite.

Transition Probabilities: For $n \leq N$

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

1.8.2 Optimal Stopping 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 decide to park in it or proceed forward. The probability that any spot is unoccupied equals p , independent of all other spots.

The optimal stopping formulation follows. Elements of S' consist of a vector, (s, k) , where s indicates the location of the parking spot and k indicates whether the spot is free (1) or not (0). We represent the set of potential parking spots by the set of integers, with 0 denoting the location of the restaurant, positive numbers denoting locations before the restaurant and negative numbers denoting locations past the restaurant. Thus, $S' = \mathbb{Z} \times \{0, 1\}$. States of the form $(s, 1)$ are the only ones where the action Q corresponding to stopping is available. No rewards are accrued if the individual does not park ($f_n((s, k)) = 0$). When the individual parks, the reward is equal to the distance from the restaurant, so $g_n((s, 1)) = -|s|$. The transition probabilities of the underlying Markov chain are given by

$$b_n((s', k')|(s, k)) = \begin{cases} p & s' = s - 1, k = 1 \\ 1 - p & s' = s - 1, k = 0. \end{cases} \quad (1.12)$$

The above description represents the problem as an infinite horizon problem. It reduces to a finite horizon problem by assuming that:

- The driver will not start looking for a spot until reaching a distance M before the restaurant, and
- Observing that if the driver reaches the restaurant and has not yet found a parking spot, the driver will most certainly take the next free one.

As a result of the second observation, we could consider the state $(0, k)$ as the terminal state, and the epoch when it is reached would correspond to the terminal decision epoch. If at location 0 there is a free spot ($k = 1$), then the driver will park and get the maximal reward of 0. If $k = 0$, the driver will continue. Since the location of the next free spot follows a geometric distribution with parameter p , the expected distance from the restaurant to the eventual parking spot will be $1/p$. Thus, the terminal reward is

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

Online Dating (also known as the Secretary Problem)

This example provides a modern take on a classical problem historically known as the *secretary problem*. An individual is searching for dates 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 swipes left ($f_n(s) = 0$). Rewards correspond to the probability of choosing the best candidate and are only received 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, let us write out the probability that $g_n(1)$ represents explicitly:

$$\begin{aligned}
 g_n(1) &:= P(\text{profile } n \text{ is rank 1} \mid \text{profile } n \text{ has the highest rank of first } n \text{ profiles}) \\
 &= \frac{P(\text{profile } n \text{ is rank 1} \cap \text{profile } n \text{ has the highest rank of first } n \text{ profiles})}{P(\text{profile } n \text{ has the highest rank of first } n \text{ profiles})} \\
 &= \frac{P(\text{profile } n \text{ is rank 1})}{P(\text{profile } n \text{ has the highest rank of first } n \text{ profiles})} \\
 &= \frac{1/N}{1/n} = \frac{n}{N}.
 \end{aligned} \tag{1.13}$$

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, we 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

$$b_n(j|s) = \begin{cases} \frac{1}{n+1} & j = 1, s \in \{0, 1\} \\ \frac{n}{n+1} & j = 0, s \in \{0, 1\} \end{cases} \quad (1.14)$$

1.9 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 ??) 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.

1.9.1 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 strategy is typically used when a team is behind by one or two goals late in the game in the hopes of tying the score and sending the game to overtime. 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 strategy is often employed late in a game 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, we describe a discrete time version. 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.

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 $w_A + w_B < 1$. We 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,

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

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: We assume the state keeps track of 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

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

We include the additional state $G + 1$ to ensure a finite state space. In our formulation, we regard this state as an absorbing state. 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 goal (action a_0). We only consider pulling the goalie (action a_1) when the goal differential is between 1 and G .

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

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: The transitions probabilities are

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

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 $p_A = p_B = 0.003125$, $w_A = 0.008875$ and $w_B = 0.02661$. On average, teams pull their goalie around 4 minutes 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 penalized 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 player in the penalty box increases Team A’s scoring rate to approximately 12 goals per 60 minutes, and decreases Team B’s scoring rate to approximately 11 goals per 60 minutes. Note that this latter quantity is still much higher than Team B’s goal scoring rate when Team A doesn’t pull its goalie, which is 6.54 goals per 60 minutes.

1.9.2 A Handicap 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, we briefly review relevant aspects of tennis scoring. 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 handicap point 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.

We now formalize the problem. 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 probability that Player A wins the point on serve (return) if it is played out. We assume the budget is for the entire match and that Player A can use a credit regardless of which player 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 is serving (Player B is serving). We include two absorbing states, W and L , that correspond to Player A winning and losing the match, respectively. Thus, the state space is

$$S = (Q \times \{0, 1, \dots, B\} \times \{0, 1\}) \cup \{W\} \cup \{L\}. \quad (1.15)$$

Each state q represents a set score, a game score and a within-game score. Note that B is upper bounded by 24 times the number of sets to win the match. In a best two out of three sets match, B is bounded by 48 and in a best three out of five sets match B is bounded by 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 .

$$A_s = \begin{cases} \{a_0, a_1\}, & s = (q, b, k) \text{ with } b > 0 \\ \{a_0\}, & \text{otherwise.} \end{cases} \quad (1.16)$$

Rewards: Since Player A's objective is 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

$$r(s, a, s') = \begin{cases} 1, & \text{if } s \neq W, a \in A_s, s' = W \\ 0, & \text{otherwise.} \end{cases} \quad (1.17)$$

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 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 $s = (q, b, k)$,

$$p(j|s, a) = \begin{cases} 1, & \text{if } j = (q^+, b - 1, k_q^+), a = a_1 \\ p_k, & \text{if } j = (q^+, b, k_q^+), a = a_0 \\ 1 - p_k, & \text{if } j = (q^-, b, k_q^-), a = a_0 \\ 0, & \text{otherwise.} \end{cases} \quad (1.18)$$

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 include 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.

1.10 The Art of Modeling

One learns to formulate Markov decision processes by studying how others have done so and practicing themselves. Formulations that appear particularly crisp are likely the result of numerous iterations in 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.

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. In our examples below, we present problem statements that are self-contained, with the information needed to fully formulate the problem.
- *Draw a timeline of events.* Carefully determine *when* the **state** information becomes available, **actions** are chosen (i.e., the **decision epochs**), **rewards** are received and **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 1.4), liver transplant model (Section 1.5), the Grid World model (Section 1.7) and the optimal stopping models (Section 1.8). 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 1.8.2).
- *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 1.6) and on-line dating application (Section 1.8.2) 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. Be sure not to include actions as part of the state space. 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 states 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.
- *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. When the reward does not depend on the subsequent state, we will leave it out of the notation.

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.

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 total 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 our Markov decision process formulation seeks to maximize rewards, we regard costs 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 1.5 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 1.2), revenue management (Section 1.1) and queuing control (Section 1.3) 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 1.4), clinical decision making (Section 1.5) and sports strategy (Section 1.9) 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 1.4), optimal parking and online dating (Section 1.8.2), pulling the goalie (Section 1.9.1), and tennis handicapping (Section 1.9.2) the reward is implicit in the chosen model objective. In other cases, rewards may be derived in consultation with the decision maker.

1.11 Bibliographic Remarks

Inventory models (Section 1.2) date back to at least Arrow et al. [1951] and Dvoretzky et al. [1952]. The book of Arrow et al. [1958] is an important early reference and Porteus [2002] provides an overview in book form. 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 model in Section 1.1 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 1.3 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 van Roy [2003] where they use it to illustrate approximate dynamic programming methods.

The lion hunting example in Section 1.4 is adopted from Clark [1987]. That paper provides many of the parameter values described in the Application Challenges section, 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 and 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 we described 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.

Our formulation of the advance scheduling model in Section 1.6 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 Goggun 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 is given in Sauré et al. [2012] in the context of radiation therapy.

A Grid World model appears in ?. 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 1.9.1 originates with Morrison [1976]. A dynamic programming formulation was provided in Washburn [1991]. Hall [2020] provides the recent data quoted in the Applications Challenge section. The tennis handicapping model in Section 1.9.2 appears in Chan and Singal [2016]. The amazing book by Kemeny and Snell [1960] includes a Markov Chain model for a tennis game.

1.12 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 1.1 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, there is an extra cost incurred.
4. Formulate a service rate control queuing model with a fixed cost C for changing the service rate. Note it is not sufficient to just modify the reward function.
5. Formulate a combined admission and service rate queuing control problem as a Markov decision process.
6. 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 $p_E + p_F < 1$. We 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.
7. (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 at queue 1 also require service at queue 2. Type 2 jobs arrive directly at queue 2 and require service at queue 2 only. **(June 22 insert picture)** In each period, no job arrives or either a type 1 job arrives, a type 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.
8. Formulate a version of the Grid World Navigation model in which there is a positive probability that the robot drops the coffee cup, which increases if the cup is full.
- (a) Suppose the cup is breakable and if it is dropped, the robot needs to return to the office to retrieve another one.
 - (b) Suppose instead that the cup is not breakable and the robot spends one epoch picking up the now empty cup.

Clearly state any assumptions you are making in formulating this model.

9. Formulate the following single machine 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 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.

10. [Stengos and Thomas, 1980] Consider the following generalization of the previous problem. You own two pieces of equipment which sometime 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?

11. Reformulate Exercise 10 assuming that the time it takes to complete maintenance on a piece on equipment is random. In particular, assume that maintenance is completed in any period with probability q , independent of other periods.
12. [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.

13. Formulate a finite horizon Markov decision process as an infinite horizon model by augmenting the state with the decision epoch and adding absorbing states.
14. In the organ transplantation problem, consider an extension where outcomes depend on the quality of the organ that is offered. Modify the formulation to account for this possibility.
15. Formulate the advance scheduling problem when appointments not booked on first day available are added to the next days demand with cost C' .

16. 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.
17. Formulate the advance scheduling problem where rewards are accrued if patients are scheduled before their fixed target date.
 - (a) Let d_k be the reward for each patient of class k who is scheduled before the target date T_k .
 - (b) What is the interpretation of the objective if $d_k = 1$ for all urgency classes k ?
18. Formulate a finite horizon version of the Grid World 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 ?
19. 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?
20. 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. We 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.
21. 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.
22. 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.

23. 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.
24. 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.
25. 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.
26. "Scrabble, like life, is a trade-off between today and tomorrow-between sending and saving. It's what an economist would call a dynamic programming problem." (Seven Games: A Human History" Oliver Roeder, W.w. Norton Company New York 2022) The game of Scrabble provides an opportunity for applying Markov decision processes. Provide 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, which we do not believe has appeared in the literature.
27. Formulate the reward function of the inventory management example when the revenue of the inventory sold is included.

(re-order the questions, but wait until solutions for chap 3 are done)

Bibliography

- O. Alagoz, L. M. Maillart, A. J. Schaefer, and M. S. Roberts. Determining the acceptance of cadaveric livers using an implicit model of the waiting list. *Operations Research*, 55:24–36, 2007.
- K. J. Arrow, D. Blackwell, and M. A. Girshick. Bayes and minimax solutions of sequential decision problems. *Econometrica*, 17:213–244, 1949.
- K. J. Arrow, T. Harris, and J. Marschak. Optimal inventory policy. *Econometrica*, 19:250–272, 1951.
- K. J. Arrow, S. Karlin, and H. E. Scarf. *Studies in the Mathematical Theory of Inventory and Production*. Stanford University Press, Stanford, CA, 1958.
- D. Bertsimas and R. Shioda. Restaurant revenue management. *Operations Research*, 51:472–486, 2003.
- V. Carter and R. E. Machol. Technical note — operations research on football. *Operations Research*, 19:541–544, 1971.
- A. Cayley. Mathematical questions with their solutions, no. 4528. *Educational Times*, 23:18, 1875.
- T. C. Y. Chan and R. Singal. A Markov Decision Process-based handicap system for tennis. *Journal of Quantitative Analysis in Sports*, 12:179–189, 2016.
- T. C. Y. Chan, C. Fernandes, and M. L. Puterman. Points gained in football: Using Markov process-based value functions to assess team performance. *Operations Research*, 69:877–894, 2021.
- Y. S. Chow, H. Robbins, and D. Siegmund. *Great Expectations: The theory of optimal stopping*. Houghton-Mifflin, New York, 1971.
- C. W. Clark. The lazy adaptable lions: a Markovian model of group foraging. *Animal Behavior*, 35:361–368, 1987.
- D. P. de Farias and B. van Roy. The linear programming approach to approximate dynamic programming. *Operations Research*, 51:850–865, 2003.

- A. Dvoretzky, J. Kiefer, and J. Wolfowitz. The inventory problem: I. case of known distributions of demand. *Econometrica*, 20:187, 1952.
- G. Gallego and G. van Ryzin. Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Management Science*, 40:999–1020, 1994.
- Y. Gocgun and M. L. Puterman. Dynamic scheduling with due dates and time windows: An application to chemotherapy patient appointment booking. *Health Care Management Science*, 17:60–76, 2014.
- M. Hall. The state of goalie pulling in the NHL. <https://hockey-graphs.com/2020/05/18/the-state-of-goalie-pulling-in-the-nhl/>, May 2020. [accessed 14-September-2021].
- D. P. Heyman. Optimal operating policies for M/G/1 queueing systems. *Operations Research*, 16:362–382, 1968.
- R. Howard. *Dynamic Programming and Markov Processes*. MIT Press, Cambridge, MA, 1960.
- T. Hubel and N. Jordan O. Dewhirst J. McNutt A. Wilson et al., J. Myatt. Energy cost and return in african wild dogs and cheetahs. *Nature Communications*, 7:11034, 2016.
- S. Karlin. Stochastic models and optimal policies for selling an asset. In K. J. Arrow, S. Karlin, and H. Scarf, editors, *Studies in Applied Probability and Management Science*, pages pp. 148–158. Stanford University Press, Palo Alto, CA, 1962.
- E. J. Kelly and P. L. Kennedy. A dynamic stochastic model of mate desertion. *Ecology*, 74:351–366, 1993.
- J. G. Kemeny and J. L. Snell. *Finite Markov Chains*. Van Nostrand-Reinhold, New York, 1960.
- M. Kurt, B. T. Denton, A. J. Schaefer, N. D. Shah, and S. A. Smith. The structure of optimal statin initiation policies for patients with type 2 diabetes. *IIE Transactions on Healthcare Systems Engineering*, 1:49–65, 2011.
- M. Mangel and C. W. Clark. Towards a unified foraging theory. *Ecology*, 67:1127–1138, 1986.
- D. G. Morrison. On the optimal time to pull the goalie: A poisson model applied to a common strategy used in ice hockey. *TIMS Studies in Management Science*, 4: 67–78, 1976.
- P. Naor. On the regulation of queue size by levying tolls. *Econometrica*, 37:15–24, 1969.

- J. Patrick, M. L. Puterman, and M. Queyranne. Dynamic multi-priority patient scheduling. *Operations Research*, 56:1507–152, 2008.
- E. Porteus. *Foundations of Stochastic Inventory Theory*. Stanford Business Books, 2002.
- M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons., 1994.
- A. Sauré, J. Patrick, and M. L. Puterman. Optimal multi-appointment scheduling. *European Journal of Operational Research*, 223:573–584, 2012.
- S. M. Shechter, M. D. Bailey, A. J. Schaefer, and M. S. Roberts. The optimal time to initiate HIV therapy under ordered health states. *Operations Research*, 56:20–33, 2008.
- E. Sirot and C. Bernstein. Time sharing between host searching and food searching in parasitoids: state-dependent optimal strategies. *Behavioral Ecology*, 7:189–194, 1996.
- B. C. Smith, J. F. Leimkuhler, and R. M. Darrow. Yield management at american airlines. *Interfaces*, 22:8–31, 1992.
- G.L. Smuts. Diet of lions and spotted hyenas assessed from stomach contents. *S. Afr. J. Wildl. Res.*, 9:19–25, 1979.
- M. J. Sobel. Optimal average-cost policy for a queue with start-up and shut-down costs. *Operations Research*, 17:145–162, 1969.
- D. Stengos and L. C. Thomas. The blast furnaces problem. *European Journal of Operational Research*, 4:330–336, 1980.
- K. Talluri and G. van Ryzin. *The Theory and Practice of Revenue Management*. Springer US, 2004.
- A. Wald. *Sequential Analysis*. John Wiley & Sons, New York, 1947.
- A. Wald and J. Wolfowitz. Optimal character of the sequential probability ratio test. *Ann. of Math. Stat.*, 19:326–339, 1948.
- A. Washburn. Still more on pulling the goalie. *Interfaces*, 21:59–64, 1991.
- M. Yadin and P. Naor. On queueing systems with variable service capacities. *Naval Research Logistics Quarterly*, 14:43–53, 1967.

Chapter 2

Index

1. xx