

Dynamic Programming

Operations Research

Anthony Papavasiliou

- 1 Multi-Stage Decision Making under Uncertainty
- 2 Dynamic Programming
- 3 Why Is Dynamic Programming Any Good?
- 4 Examples
 - The Monty Hall Problem
 - Pricing Financial Securities

- 1 Multi-Stage Decision Making under Uncertainty
- 2 Dynamic Programming
- 3 Why Is Dynamic Programming Any Good?
- 4 Examples
 - The Monty Hall Problem
 - Pricing Financial Securities

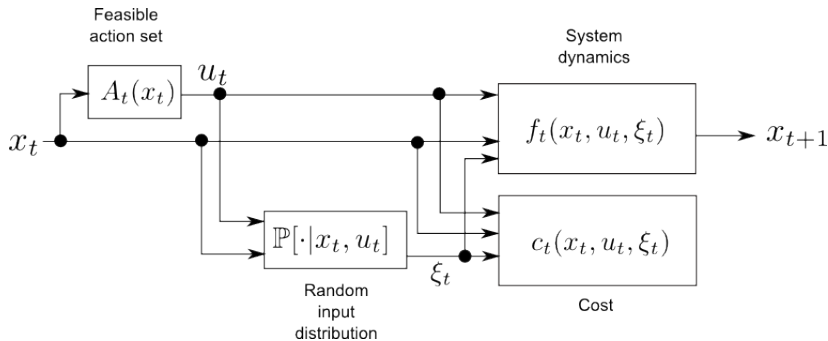
- **Dynamical** system with $H < \infty$ discrete time stages
 - Extensions exist for infinite horizon ($H = \infty$)
 - Extensions exist for continuous time
- **Controlled** system, denote u_t as *continuous* decision at stage t
- **Stochastic** system, denote ξ_t as *discrete* random vector at stage t
 - Extensions exist for continuous uncertainty
- Denote x_t as *continuous state* of the system at the *end* of stage t
 - State encodes everything we need to know, except ξ_t and u_t , for describing the evolution of the system
- Transition equation:

$$x_{t+1} = f_t(x_t, u_t, \xi_t)$$

Setting (II)

- *Markovian* uncertainty: we can define probability distribution $\mathbb{P}[\cdot|x_t, u_t]$ for ξ_t , *independently* of $\xi_{t-1}, \xi_{t-2}, \dots, \xi_0$
- Denote $A_t(x_t)$ as set of finite actions at stage t
- Costs are *additive*, denote $c_t(x_t, u_t, \xi_t)$ as cost per time stage
- Usually (but not always), we will assume that x_t, ξ_t , and u_t live in \mathbb{R}^n

Block Diagram Representation



- The flow of information is consistent (everything depends on information that is already revealed)
- The process is repeated identically over stages

Sequence of Events

In a given time stage,

- 1 observe state x_t
- 2 decide u_t *after* observing x_t
- 3 sample ξ_t from a distribution that depends on x_t, u_t
- 4 Incur cost $c_t(x_t, u_t, \xi_t)$
- 5 Move to new state x_{t+1}

This is **not** your ‘usual’ optimization

- In ‘usual’ optimization we are looking for an optimal *vector* x^*
- In multi-stage optimization under uncertainty we are looking for a *sequence of functions* $\mu_t(x_t)$
- The functions $\mu_t(x_t)$ are called a **policy**, they tell us what to do if we observe x_t in stage t

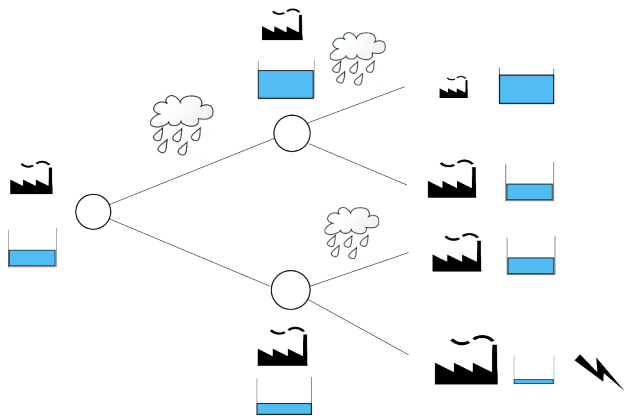
- Recall costs are additive
 - For $t = 0, \dots, H - 1$, we incur cost $c_t(x_t, u_t, \xi_t)$
 - Assume final-period cost only depends on x_H , i.e. $c_H(x_H)$
- **Key observation:** *given* a policy $\mu_t(x_t)$, we can define a distribution for the sequence $(x_t, \xi_t), t = 0, \dots, H$
- Given a distribution for the sequence $(x_t, \xi_t), t = 0, \dots, H$, we can define expected cost

$$\mathbb{E}\left[\sum_{t=0}^{H-1} c_t(x_t, \mu_t(x_t), \xi_t) + c_H(x_H)\right]$$

We are looking for the **optimal policy**: the *policy* which minimizes *expected* cost

$$(MP) : \min_{\mu_t} \mathbb{E} \left[\sum_{t=0}^{H-1} c_t(x_t, \mu_t(x_t), \xi_t) + c_H(x_H) \right]$$
$$\mu_t(x_t) \in \mathbf{A}_t(x_t)$$

Recalling Hydrothermal Scheduling



- Too much water in dams leads to water spillage and unnecessary thermal generation costs
- Too little water in dams leads to load curtailment

Hydrothermal Scheduling Problem Statement

- Time-varying electricity demand D_t
- Three options
 - Hydro units: produce q_t at zero cost
 - Thermal units: produce p_t at marginal cost C
 - Load shedding: cut supply by l_t at marginal cost V
- Rainfall uncertainty: independent identical normal distribution with mean μ and standard deviation σ
- Hydro reservoir can store up to E units of energy
- Thermal generators can produce up to P units of power per period

Hydrothermal Scheduling Model Description

- *Continuous* action vector: $u_t = (p_t, q_t, l_t) \in \mathbb{R}^3$
- *Continuous* state vector $x_t \in \mathbb{R}$: level of reservoir at the *beginning* of stage t
- Feasible action set:

$$A_t(x_t) = \{(p_t, q_t, l_t) :$$

$$q_t \leq x_t$$

$$p_t + q_t + l_t = D_t$$

$$p_t \leq P$$

$$p_t, q_t, l_t \geq 0\}$$

- *Continuous* random disturbance $\xi_t \in \mathbb{R}$: rainfall

- Transition probability function:

$$\mathbb{P}[\xi_t \leq R] = \Phi\left(\frac{R - \mu}{\sigma}\right),$$

where $\Phi(\cdot)$ is the cdf of a standard normal random variable

- System transition function:

$$x_{t+1} = f_t(x_t, u_t, \xi_t) = \min(E, x_t + \xi_t - q_t)$$

- Cost function:

$$c_t(x_t, u_t, \xi_t) = C \cdot p_t + V \cdot I_t$$

Hydrothermal Scheduling with AR Rainfall

Same problem as before, except rainfall r_t follows an autoregressive (AR) process:

$$r_t = c + \phi \cdot r_{t-1} + w_t$$

- c and ϕ are fixed parameters
- w_t : independent identical distribution according to a normal distribution with mean μ and standard deviation σ

Hydrothermal Scheduling with AR Rainfall: Model Formulation

- Redefine random disturbance as $\xi_t = w_t \in \mathbb{R}$
- State of the system: $x_t = (e_t, r_t)^T \in \mathbb{R}^2$
 - e_t : level of energy stored in the reservoir
 - r_t : rainfall
- System dynamic function:

$$x_{t+1} = (e_{t+1}, r_{t+1})^T = f_t(x_t, u_t, \xi_t) = \begin{bmatrix} \min(E, x_t + \xi_t - q_t) \\ c + \phi \cdot r_t + \xi_t \end{bmatrix}$$

Capacity Expansion Problem

- *Continuous* action vector $u_t = (z_t, y_t) \in \mathbb{R}^{nm+n-1}$
 - Amount of capacity constructed:
 $z_{it} = (z_{1t}, \dots, z_{n-1,t}) \in \mathbb{R}^{n-1}$
 - Amount of power that from technology i to block j :
 $y_t = (y_{11t}, \dots, y_{1mt}, \dots, y_{n1t}, \dots, y_{nmt}) \in \mathbb{R}^{nm}$
- *Continuous* state vector: capacity that has been constructed so far for each technology,
 $x_t = v_t = (v_{1t}, \dots, v_{n-1,t}) \in \mathbb{R}^{n-1}$.
- *Discrete* uncertain demand $D_t = (D_{1t}, \dots, D_{mt}) \in \mathbb{R}^m$ with distribution $\mathbb{P}[\cdot]$, independent of x_t and u_t

Capacity Expansion Problem (II)

- Feasible action set:

$$\begin{aligned} A_t(x_t) = \{ & (z_t, y_t) : \\ & \sum_{j=1}^m y_{ijt} \leq x_{it}, i = 1, \dots, n-1, \\ & \sum_{i=1}^n y_{ijt} = D_j, j = 1, \dots, m \\ & y_t \geq 0, z_t \geq 0 \} \end{aligned}$$

- System transition function:

$$x_{i,t+1} = x_{it} + z_{it}, i = 1, \dots, n-1$$

Capacity Expansion Problem (III)

- Cost function:

$$C_t(x_t, u_t, \xi_t) = \sum_{i=1}^{n-1} I_i \cdot z_{it} + \sum_{i=1}^n \sum_{j=1}^m C_i \cdot T_j \cdot y_{ijt},$$

where

- I_i : investment cost of technology i
 - C_i : marginal cost of technology i
 - T_j : (deterministic) duration of block j
- Note: capacity built in period t cannot be used for satisfying the demand of period t

Machine Scheduling: Problem Statement

- Machine produces P units of output when on
- Cost of C is paid for every period that the machine is on
- Machine output earns time-varying price λ_t
- Machine needs to stay on for at least 3 hours once started up

Machine Scheduling: Model Description

- Action set: $\mathbb{B} = \{\text{Stay}, \text{Change}\}$
- State: number of hours that have elapsed since the machine was last turned on, belongs to set $\mathbb{Z} = \{0, 1, 2, \dots\}$ - 0 belongs to 'Off'
- Feasible action set:

$$A_t(0) = \{\text{Stay}, \text{Change}\},$$

$$A_t(x_t) = \{\text{Stay}\}, x_t = 1, 2$$

$$A_t(x_t) = \{\text{Stay}, \text{Change}\}, x_t \geq 3$$

- System transition function:

$$x_{t+1} = f_t(0, \text{Stay}) = 0$$

$$x_{t+1} = f_t(x_t, \text{Stay}) = x_t + 1, x_t \geq 1$$

$$x_{t+1} = f_t(0, \text{Change}) = 1$$

$$x_{t+1} = f_t(x_t, \text{Change}) = 0, x_t \geq 1$$

- Cost function:

$$c_t(x_t, u_t) = (C - \lambda_t \cdot P), x_t \geq 1$$

$$c_t(0, u_t) = 0$$

Table of Contents

- 1 Multi-Stage Decision Making under Uncertainty
- 2 Dynamic Programming**
- 3 Why Is Dynamic Programming Any Good?
- 4 Examples
 - The Monty Hall Problem
 - Pricing Financial Securities

- Solving (*MP*) means solving for a policy / mapping μ_t , not a vector u_t
- **Value function** $V_t(x_t)$: least expected cost if optimal decisions would be made from stage t onwards *given* state x_t

The Dynamic Programming Algorithm

Dynamic programming algorithm:

- Starting from $t = H$, for all $x_t \in A_t(x_t)$, compute

$$V_H(x_H) = c_H(x_H).$$

- Moving backwards in time, for all $t = H - 1, \dots, 0$, for all $x_t \in A_t(x_t)$, compute

$$V_t(x_t) = \min_{u_t \in A_t(x_t)} \mathbb{E}_{\xi_t} [(c_t(x_t, u_t, \xi_t) + V_{t+1}(f_{t+1}(x_t, u_t, \xi_t))) | x_t, u_t]$$

where the expectation is over the distribution of ξ_t given u_t and x_t

Intuition: an optimal policy over a horizon $\{0, \dots, H\}$ is optimal for $\{t, \dots, H\}$

Value functions $V_t(x_t)$ allow decomposition of multi-period problem to single-stage optimization problems

Define **Q functions**:

$$\begin{aligned} Q_t(x_t, \xi_t) &= \min_{u_t, x_{t+1}} c_t(x_t, u_t, \xi_t) + V_{t+1}(x_{t+1}) \\ &\text{s.t. } x_{t+1} = f_t(x_t, u_t, \xi_t) \\ &\quad u_t \in A_t(x_t) \end{aligned}$$

Interpretation of $Q_{t+1}(x_t, \xi_t)$: cost of being in x_t given that ξ_t *will occur* in t

Q function can be used to solve problem backwards:

$$V_t(x_t) = \mathbb{E}_{\xi_t}[Q_t(x_t, \xi_t)|x_t, u_t].$$

Curse of Dimensionality

Consider discretization of each component of $x_t \in \mathbb{R}^m$, $u_t \in \mathbb{R}^n$, $\xi_t \in \mathbb{R}^p$ into d points

At stage t , computation of $V_t(x_t)$ for all x_t requires

- for all d^m possible values of x_t
- compute expectation \Rightarrow summation over d^p values of ξ_t
- minimization \Rightarrow comparison of d^n possible values of u_t

Each stage of DP algorithm requires $O(d^{m+n+p})$ operations \Rightarrow overall complexity of $O(H \cdot d^{m+n+p})$

Table of Contents

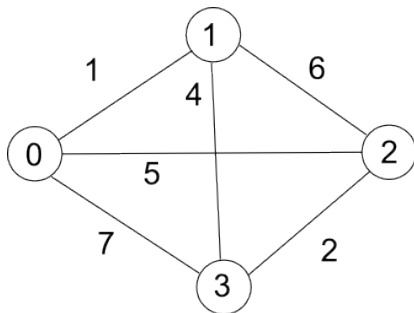
- 1 Multi-Stage Decision Making under Uncertainty
- 2 Dynamic Programming
- 3 Why Is Dynamic Programming Any Good?**
- 4 Examples
 - The Monty Hall Problem
 - Pricing Financial Securities

Recall that central entity of DP algorithm is the *value function*

Main idea of DP efficiency: avoid unnecessary repetition of computation by storing future cost data in value functions

Traveling Salesman Problem

Goal: starting from city 0, find *minimum distance* tour that goes through all cities *exactly once* and returns to 0



c_{ij} : distance from city i to city j (indicated on arcs)

Tour	Distance
01320	12
02310	12
01230	16
03210	16
02130	22
03120	22

Given $H + 1$ cities:

- must examine $H!$ tours
- Computation of cost of each tour: H summations

Complexity of enumeration: $O(H! \cdot H)$

Idea: interpret each city as one 'stage' in multi-stage decision

State: information necessary for deciding next move

- R_t : set of cities that still need to be visited
- i : current city

Value function $V_t(R_t, i)$: most efficient way of visiting cities in R_t exactly once, starting from i and ending in 0

$$V_t(R_t, i) = \min_{j \in R_t} c_{ij} + V_{t+1}(R_t - \{j\}, j)$$

Note: $V_t(R_t, i)$ is reused

Complexity of Dynamic Programming for TSP

At stage t , computation of V_t for all i , R_t requires:

- for H different values of i
- for $\binom{H}{H-t}$ different values of R_t
- one minimization over $H-t$ values (size of R_t)

Total number of operations:

- For $V_0(\{1, \dots, H\}, 0)$: H summations
- For $1 \leq t \leq H$:

$$\sum_{t=1}^H H \cdot \binom{H}{H-t} \cdot (H-t+1) = O(H^2 \cdot 2^H)$$

Note: Complexity remains exponential, better than factorial

Where did the computational savings come from?

Value function	Evaluation
$V_3(\emptyset, 1)$	1
$V_3(\emptyset, 2)$	5
$V_3(\emptyset, 3)$	7
$V_2(\{2\}, 1) = c_{12} + V_3(\emptyset, 2)$	11
$V_2(\{3\}, 1)$	11
$V_2(\{1\}, 2)$	7
$V_2(\{3\}, 2)$	9
$V_2(\{1\}, 3)$	5
$V_2(\{2\}, 3)$	7
$V_1(\{2, 3\}, 1) = \min\{c_{12} + V_2(\{3\}, 2), c_{13} + V_2(\{2\}, 3)\}$	11
$V_1(\{1, 3\}, 2)$	7
$V_1(\{1, 2\}, 3)$	9
$V_0(\{1, 2, 3\}, 0) = \min_{x \in \{1, 2, 3\}} (c_{0x} + V_2(\{1, 2, 3\} - \{x\}, x))$	12

Table of Contents

- 1 Multi-Stage Decision Making under Uncertainty
- 2 Dynamic Programming
- 3 Why Is Dynamic Programming Any Good?
- 4 **Examples**
 - The Monty Hall Problem
 - Pricing Financial Securities

The Monty Hall Problem

- 1 the player is asked to pick a curtain
- 2 the host opens up a curtain with a goat behind it
- 3 the player can keep the curtain that she chose originally, or switch to the remaining curtain
- 4 the player keeps the content behind the curtain that was selected in step (iii)

Should the player change curtains in step (iii), or not?

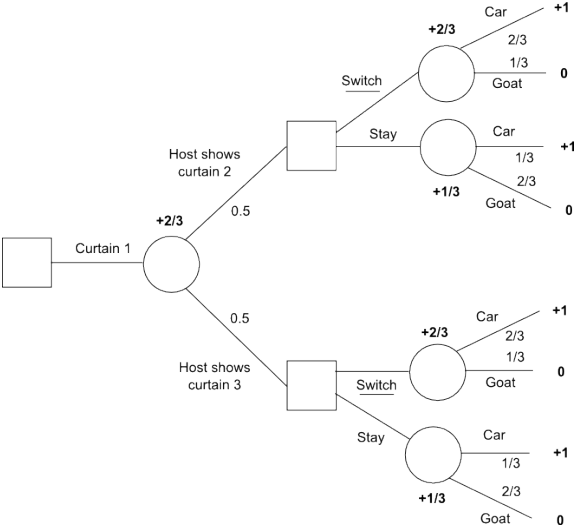
Solution of the Monty Hall Problem

Assumption: when the host opens a curtain in step (ii), the host will choose curtains with equal likelihood if both of the curtains not chosen by the player hide a goat

Thinking about the decision tree:

- We need three stages
- Symmetry \Rightarrow we do not lose generality by assuming that the player picks curtain 1 in the first step
- Symmetry + assumption \Rightarrow equal probability of host opening curtain 2 or curtain 3
- Uncertainty in stage 1 is *not* the location of the sports car, this cannot be observed!
- Compute transition probabilities of second stage using Bayes' theorem

Decision Tree of the Newsboy Problem



Probability of Winning if We Stay

Bayes' theorem:

$$\begin{aligned}\mathbb{P}[\text{Car in C1} | \text{Host shows C2}] &= \\ \frac{\mathbb{P}[\text{Car in C1, Host shows C2}]}{\mathbb{P}[\text{Host shows C2}]} &= \\ \frac{\mathbb{P}[\text{Host shows C2} | \text{Car in C1}] \cdot \mathbb{P}[\text{Car in C1}]}{\mathbb{P}[\text{Host shows C2}]} &= \\ \frac{1/2 \cdot 1/3}{1/2} &= \frac{1}{3}\end{aligned}$$

Probability of winning if we stay = original probability of winning
 \Rightarrow we have not gained (or lost) anything by staying

Intuitive? Maybe ...

Probability of Winning if We Switch

Bayes' theorem:

$$\begin{aligned}\mathbb{P}[\text{Car in C3} | \text{Host shows C2}] &= \\ \frac{\mathbb{P}[\text{Car in C3, Host shows C2}]}{\mathbb{P}[\text{Host shows C2}]} &= \\ \frac{\mathbb{P}[\text{Host shows C2} | \text{Car in C3}] \cdot \mathbb{P}[\text{Car in C3}]}{\mathbb{P}[\text{Host shows C2}]} &= \\ \frac{1 \cdot 1/3}{2/3} &= \frac{2}{3}\end{aligned}$$

Chances of winning double if we switch?

Intuitive? Try this: the host deliberately leaves one door unrevealed

Someone might argue that switching gives a 50/50 chance of winning, because it is like picking from two doors. This intuition is wrong!

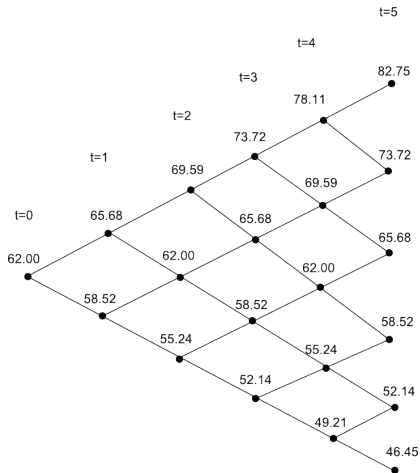
American call option: financial instrument that allows its owner to buy a certain financial asset at a *strike price* at or before a certain *expiration date*

Call option at time t is worth $\max(S_t - k, 0)$, where

- S_t : price of financial asset at time t
- k : strike price of the option

Use DP in order to determine how much an option is worth at time $t = 0$

Lattice Model of Stock Price S_t



Upward transition probability is $q = 0.5577$, downward transition probability is $1 - q$

Backward Solution

Denote $V_t(i)$ as the value of the stock at stage t , and state i , where i corresponds to one of the nodes in the lattice

Consider strike price $k = 60$

Period 5 payoff:

$$V_5(1) = 82.75 - 60 = 22.75$$

$$V_5(2) = 73.72 - 60 = 13.72$$

$$V_5(3) = 65.68 - 60 = 5.68$$

$$V_5(4) = 0$$

$$V_5(5) = 0$$

$$V_5(6) = 0$$

Backward Solution (II)

Period 4 payoff:

$$V_4(1) = \max(78.11 - 60, \mathbb{E}[V_5(j)|i = 1]) = 18.7560$$

$$V_4(2) = \max(69.59 - 60, \mathbb{E}[V_5(j)|i = 2]) = 10.1639$$

$$V_4(3) = \max(62 - 60, \mathbb{E}[V_5(j)|i = 3]) = 3.1677$$

$$V_4(4) = 0$$

$$V_4(5) = 0$$

Period 3 payoff:

$$V_3(1) = \max(73.72 - 60, \mathbb{E}[V_4(j)|i = 1]) = 14.9557$$

$$V_3(2) = \max(65.68 - 60, \mathbb{E}[V_4(j)|i = 2]) = 7.0695$$

$$V_3(3) = \max(0, \mathbb{E}[V_4(j)|i = 3]) = 1.7666$$

$$V_3(4) = 0$$

Backward Solution (III)

Period 2 payoff:

$$V_2(1) = \max(69.59 - 60, \mathbb{E}[V_3(j)|i = 1]) = 11.4676$$

$$V_2(2) = \max(62 - 60, \mathbb{E}[V_3(j)|i = 2]) = 4.7240$$

$$V_2(3) = 0.9852$$

Period 1 payoff:

$$V_1(1) = \max(65.68 - 60, \mathbb{E}[V_2(j)|i = 1]) = 8.4849$$

$$V_1(2) = \max(0, \mathbb{E}[V_2(j)|i = 2]) = 3.0703$$

Period 0 payoff:

$$V_0(1) = \max(62 - 60, \mathbb{E}[V_1(j)|i = 1]) = 6.09$$