

Experimental design and Markov decision processes

The following problems

- ▶ Shortest path problems.
- ▶ Optimal stopping problems.
- ▶ Reinforcement learning problems.
- ▶ Experiment design (clinical trial) problems
- ▶ Advertising.

can be all formalised as **Markov decision processes**.

Applications

- ▶ Robotics.
- ▶ Economics.
- ▶ Automatic control.
- ▶ Resource allocation

Contents

Bandit problems

Introduction

Bernoulli bandits

Markov processes

Markov processes

Markov decision processes and reinforcement learning

Markov decision processes

Value functions

Examples

Episodic problems

Policy evaluation

Backwards induction

Continuing, discounted problems

Markov chain theory for discounted problems

Infinite horizon MDP Algorithms

The n meteorologists problem

Setting

- ▶ n meteorologists.

The n meteorologists problem

Setting

- ▶ n meteorologists.
- ▶ At time t :

The n meteorologists problem

Setting

- ▶ n meteorologists.
- ▶ At time t :
- ▶ The i -th meteorologist gives a prediction $x_{t,i} \in \{\text{dry}, \text{wet}\}$ for today's weather, y_t

The n meteorologists problem

Setting

- ▶ n meteorologists.
- ▶ At time t :
- ▶ The i -th meteorologist gives a prediction $x_{t,i} \in \{\text{dry}, \text{wet}\}$ for today's weather, y_t
- ▶ You must decide whether or not to commute by bike that day $a_t \in \{\text{bike}, \text{tram}\}$

The n meteorologists problem

Setting

- ▶ n meteorologists.
- ▶ At time t :
- ▶ The i -th meteorologist gives a prediction $x_{t,i} \in \{\text{dry}, \text{wet}\}$ for today's weather, y_t
- ▶ You must decide whether or not to commute by bike that day $a_t \in \{\text{bike}, \text{tram}\}$
- ▶ You obtain a reward $r_t = 1$ if it's dry and you bike, $r_t = -1$ if it's wet and you bike, and $r_t = 0$ otherwise.

The n meteorologists problem

Setting

- ▶ n meteorologists.
- ▶ At time t :
- ▶ The i -th meteorologist gives a prediction $x_{t,i} \in \{\text{dry}, \text{wet}\}$ for today's weather, y_t
- ▶ You must decide whether or not to commute by bike that day $a_t \in \{\text{bike}, \text{tram}\}$
- ▶ You obtain a reward $r_t = 1$ if it's dry and you bike, $r_t = -1$ if it's wet and you bike, and $r_t = 0$ otherwise.
- ▶ You then store the information y_t about the weather and the meteorologists' predictions.

Utility

$$U = \sum_t r_t$$

The n meteorologists problem is simple, as:

- ▶ You always see their predictions, as well as the weather, no matter whether you bike or take the tram (full information)
- ▶ Your actions do not influence their predictions (independence events)

In the remainder, we'll see two settings where decisions are made with either **partial information** or in a **dynamical system**. Both of these settings can be formalised with Markov decision processes.

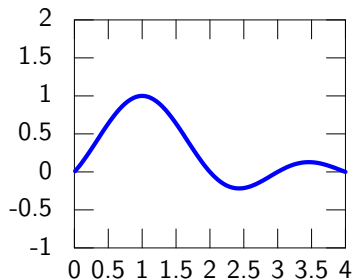
Bandit problems



Bandit problems

Applications

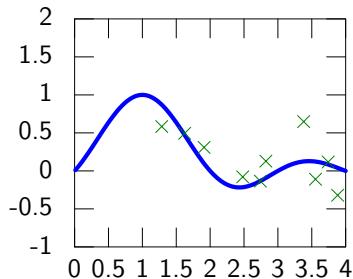
- Efficient optimisation.



Bandit problems

Applications

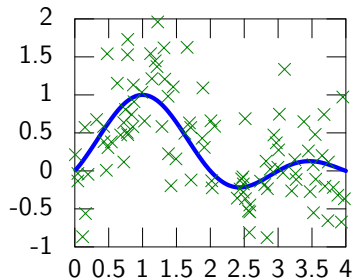
- Efficient optimisation.



Bandit problems

Applications

- Efficient optimisation.



Bandit problems

Applications

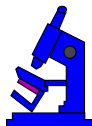
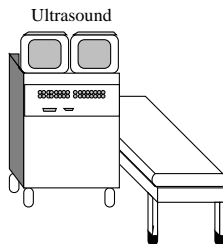
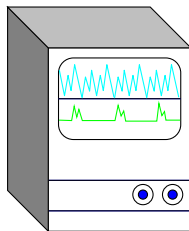
- ▶ Efficient optimisation.
- ▶ Online advertising.



Bandit problems

Applications

- ▶ Efficient optimisation.
- ▶ Online advertising.
- ▶ Clinical trials.



Bandit problems

Applications

- ▶ Efficient optimisation.
- ▶ Online advertising.
- ▶ Clinical trials.
- ▶ ROBOT SCIENTIST.



The stochastic n -armed bandit problem

Actions and rewards

- ▶ A set of **actions** $\mathcal{A} = \{1, \dots, n\}$.
- ▶ Each action gives you a **random reward** with distribution $\mathbb{P}(r_t \mid a_t = i)$.
- ▶ The **expected reward** of the i -th arm is $\rho_i \triangleq \mathbb{E}(r_t \mid a_t = i)$.

Interaction at time t

1. You choose an action $a_t \in \mathcal{A}$.
2. You observe a random reward r_t drawn from the i -th arm.

The utility is the **sum of the rewards** obtained

$$U \triangleq \sum_t r_t.$$

We must maximise the expected utility, **without knowing** the values ρ_i .

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history

$$h_t \triangleq a_1, r_1, \dots, a_t, r_t$$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

Exercise 1

Why should our action depend on the complete history?

- A *The next reward depends on all the actions we have taken.*
- B *We don't know which arm gives the highest reward.*
- C *The next reward depends on all the previous rewards.*
- D *The next reward depends on the complete history.*
- E *No idea.*

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history $h_t \triangleq a_1, r_1, \dots, a_t, r_t$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

Example 2 (The expected utility of a uniformly random policy)

If $\mathbb{P}^\pi(a_{t+1} \mid \cdot) = 1/n$ for all t , then

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history

$$h_t \triangleq a_1, r_1, \dots, a_t, r_t$$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

Example 2 (The expected utility of a uniformly random policy)

If $\mathbb{P}^\pi(a_{t+1} \mid \cdot) = 1/n$ for all t , then

$$\mathbb{E}^\pi U = \mathbb{E}^\pi \left(\sum_{t=1}^T r_t \right) = \sum_{t=1}^T \mathbb{E}^\pi r_t = \sum_{t=1}^T \sum_{i=1}^n \frac{1}{n} \rho_i = \frac{T}{n} \sum_{i=1}^n \rho_i$$

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history

$$h_t \triangleq a_1, r_1, \dots, a_t, r_t$$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

The expected utility of a general policy

$$\mathbb{E}^\pi U = \mathbb{E}^\pi \left(\sum_{t=1}^T r_t \right)$$

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history

$$h_t \triangleq a_1, r_1, \dots, a_t, r_t$$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

The expected utility of a general policy

$$\mathbb{E}^\pi U = \mathbb{E}^\pi \left(\sum_{t=1}^T r_t \right) = \sum_{t=1}^T \mathbb{E}^\pi(r_t) \quad (1.1)$$

Policy

Definition 1 (Policies)

A policy π is **an algorithm for taking actions** given the observed history

$$h_t \triangleq a_1, r_1, \dots, a_t, r_t$$

$$\mathbb{P}^\pi(a_{t+1} \mid h_t)$$

is the probability of the next action a_{t+1} .

The expected utility of a general policy

$$\begin{aligned} \mathbb{E}^\pi U &= \mathbb{E}^\pi \left(\sum_{t=1}^T r_t \right) = \sum_{t=1}^T \mathbb{E}^\pi(r_t) \\ &= \sum_{t=1}^T \sum_{a_t \in \mathcal{A}} \mathbb{E}(r_t \mid a_t) \sum_{h_{t-1}} \mathbb{P}^\pi(a_t \mid h_{t-1}) \mathbb{P}^\pi(h_{t-1}) \end{aligned} \tag{1.1}$$

Contents

Bandit problems

Introduction

Bernoulli bandits

Markov processes

Markov processes

Markov decision processes and reinforcement learning

Markov decision processes

Value functions

Examples

Episodic problems

Policy evaluation

Backwards induction

Continuing, discounted problems

Markov chain theory for discounted problems

Infinite horizon MDP Algorithms

Bernoulli bandits

Example 2 (Bernoulli bandits)

Consider n Bernoulli distributions with parameters ω_i ($i = 1, \dots, n$) such that $r_t \mid a_t = i \sim \text{Bern}(\omega_i)$. Then,

$$\mathbb{P}(r_t = 1 \mid a_t = i) = \omega_i \qquad \mathbb{P}(r_t = 0 \mid a_t = i) = 1 - \omega_i \qquad (1.2)$$

Then the expected reward for the i -th bandit is $\rho_i \triangleq \mathbb{E}(r_t \mid a_t = i) = ?$.

Bernoulli bandits

Example 2 (Bernoulli bandits)

Consider n Bernoulli distributions with parameters ω_i ($i = 1, \dots, n$) such that $r_t \mid a_t = i \sim \text{Bern}(\omega_i)$. Then,

$$\mathbb{P}(r_t = 1 \mid a_t = i) = \omega_i \qquad \mathbb{P}(r_t = 0 \mid a_t = i) = 1 - \omega_i \qquad (1.2)$$

Then the expected reward for the i -th bandit is $\rho_i \triangleq \mathbb{E}(r_t \mid a_t = i) = \omega_i$.

Bernoulli bandits

Example 2 (Bernoulli bandits)

Consider n Bernoulli distributions with parameters ω_i ($i = 1, \dots, n$) such that $r_t \mid a_t = i \sim \text{Bern}(\omega_i)$. Then,

$$\mathbb{P}(r_t = 1 \mid a_t = i) = \omega_i \qquad \mathbb{P}(r_t = 0 \mid a_t = i) = 1 - \omega_i \qquad (1.2)$$

Then the expected reward for the i -th bandit is $\rho_i \triangleq \mathbb{E}(r_t \mid a_t = i) = \omega_i$.

Exercise 1 (The optimal policy under **perfect** knowledge)

If we know ω_i for all i , what is the best policy?

- A At every step, play the bandit i with the greatest ω_i .*
- B Prefer bandits i with larger ω_i , but play them all.*
- C It depends on the horizon T .*
- D Prefer bandits i which you have played the least so far, but play them all.*
- E It is too complicated.*

The unknown reward case

Say you keep a running average of the reward obtained by each arm

$$\hat{\rho}_{t,i} = R_{t,i}/n_{t,i}$$

where $n_{t,i}$ is the number of times you played arm i and $R_{t,i}$ the total reward received from i so that whenever you play $a_t = i$:

$$R_{t+1,i} = R_{t,i} + r_t, \quad n_{t+1,i} = n_{t,i} + 1.$$

The unknown reward case

Say you keep a running average of the reward obtained by each arm

$$\hat{p}_{t,i} = R_{t,i}/n_{t,i}$$

where $n_{t,i}$ is the number of times you played arm i and $R_{t,i}$ the total reward received from i so that whenever you play $a_t = i$:

$$R_{t+1,i} = R_{t,i} + r_t, \quad n_{t+1,i} = n_{t,i} + 1.$$

Exercise 2 (The optimal policy under imperfect knowledge)

If we just keep track of the averages $\hat{p}_{t,i}$, for all i , what is the best policy?

- A At every step, play the bandit i with the greatest $\hat{p}_{t,i}$.*
- B Prefer bandits i with larger $\hat{p}_{t,i}$, but play them all.*
- C It depends on the horizon T .*
- D Prefer bandits i with smaller $n_{t,i}$, but play them all.*
- E It is too complicated.*

The unknown reward case

Say you keep a running average of the reward obtained by each arm

$$\hat{\rho}_{t,i} = R_{t,i}/n_{t,i}$$

where $n_{t,i}$ is the number of times you played arm i and $R_{t,i}$ the total reward received from i so that whenever you play $a_t = i$:

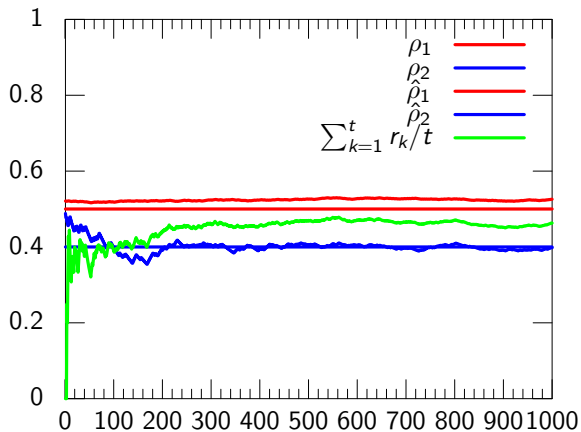
$$R_{t+1,i} = R_{t,i} + r_t, \quad n_{t+1,i} = n_{t,i} + 1.$$

You could choose to play the strategy

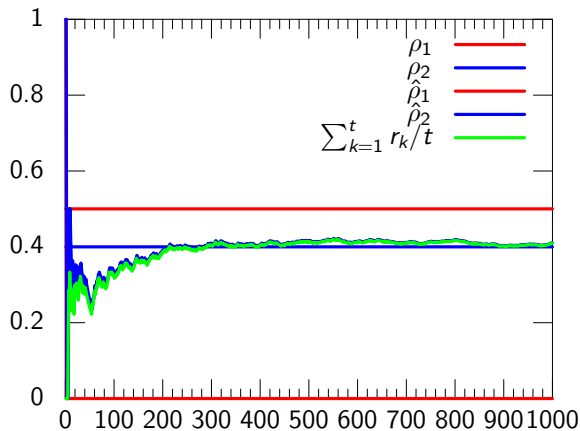
$$a_t = \arg \max_i \hat{\rho}_{t,i}.$$

where we use **non-zero** initial values $n_{0,i}, R_{0,i}$!

The uniform policy

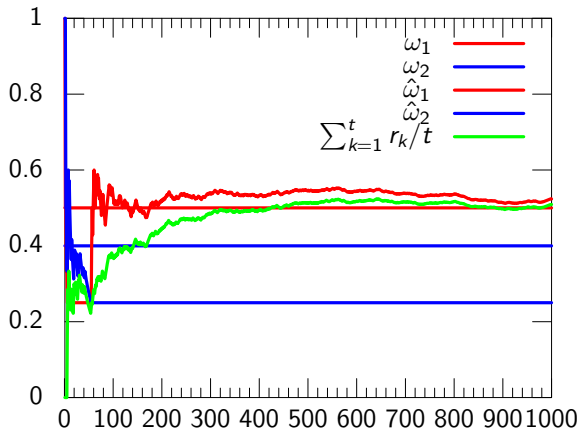


The greedy policy



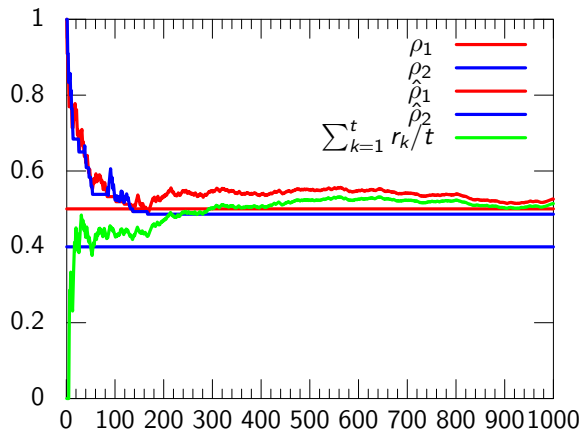
For $n_{0,i} = R_{0,i} = 0$

The greedy policy



For $n_{0,i} = R_{0,i} = 1$

The greedy policy



For $n_{0,i} = R_{0,i} = 10$

Summary

- ▶ Bandit problems are the simplest type of **partial information** problems.
- ▶ **Learning policies** for such problems must remember the complete history.
- ▶ If we know the problem parameters, simple stationary policies are optimal.
- ▶ If we don't, then our policies must carefully balance:
 - ▶ Exploration: Learning more about the problem.
 - ▶ Exploitation: Using what is already known.

From now on, we focus on the case where the problem is perfectly known.

Contents

Bandit problems

- Introduction

- Bernoulli bandits

Markov processes

- Markov processes

Markov decision processes and reinforcement learning

- Markov decision processes

- Value functions

- Examples

Episodic problems

- Policy evaluation

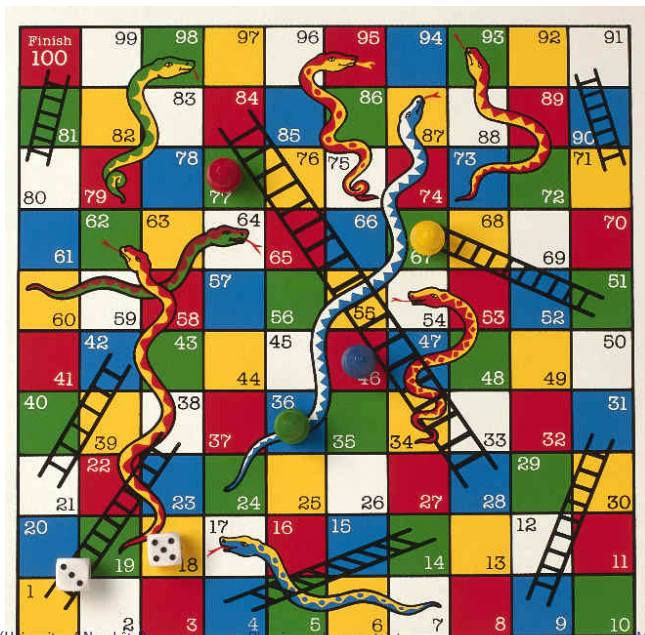
- Backwards induction

Continuing, discounted problems

- Markov chain theory for discounted problems

- Infinite horizon MDP Algorithms

A gentle introduction to Markov processes



Markov process



Definition 3 (Markov Process – or Markov Chain)

The sequence $\{s_t \mid t = 1, \dots\}$ of random variables $s_t : \Omega \rightarrow \mathcal{S}$ is a Markov process if

$$\mathbb{P}(s_{t+1} \mid s_t, \dots, s_1) = \mathbb{P}(s_{t+1} \mid s_t). \quad (2.1)$$

- ▶ s_t is **state** of the Markov process at time t .
- ▶ $\mathbb{P}(s_{t+1} \mid s_t)$ is the **transition kernel** of the process.

Markov process



Definition 3 (Markov Process – or Markov Chain)

The sequence $\{s_t \mid t = 1, \dots\}$ of random variables $s_t : \Omega \rightarrow \mathcal{S}$ is a Markov process if

$$\mathbb{P}(s_{t+1} \mid s_t, \dots, s_1) = \mathbb{P}(s_{t+1} \mid s_t). \quad (2.1)$$

- ▶ s_t is **state** of the Markov process at time t .
- ▶ $\mathbb{P}(s_{t+1} \mid s_t)$ is the **transition kernel** of the process.

Exercise 2 (Finite state machine with random input)

Let $\omega = \omega_1, \dots, \omega_t$ be an infinitely long random string of bits, $\Omega = \mathcal{S} = \{0, 1\}$ and:

$$s_{t+1} = s_t \oplus \omega_t.$$

Is s_t a Markov process?

Contents

Bandit problems

- Introduction

- Bernoulli bandits

Markov processes

- Markov processes

Markov decision processes and reinforcement learning

- Markov decision processes

- Value functions

- Examples

Episodic problems

- Policy evaluation

- Backwards induction

Continuing, discounted problems

- Markov chain theory for discounted problems

- Infinite horizon MDP Algorithms

Reinforcement learning

The reinforcement learning problem.

Learning to act in an **unknown** environment, by **interaction** and **reinforcement**.

- ▶ The environment has a changing state s_t .
- ▶ The agents observes the state s_t .
- ▶ The agent takes action a_t .
- ▶ It receives rewards r_t .

The goal (informally)

Maximise total reward $\sum_t r_t$

Types of environments

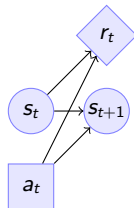
- ▶ **Markov decision processes** (MDPs).
- ▶ Partially observable MDPs (POMDPs).
- ▶ (Partially observable) **Markov games**.

Markov decision processes

Markov decision processes (MDP) μ .

At each time step t :

- ▶ We observe **state** $s_t \in \mathcal{S}$.
- ▶ We take **action** $a_t \in \mathcal{A}$.
- ▶ We receive a **reward** $r_t \in \mathbb{R}$.



Markov property of the reward and state distribution

$$\mathbb{P}_{\mu}(s_{t+1} \mid s_t, a_t)$$

$$\mathbb{P}_{\mu}(r_t \mid s_t, a_t)$$

(Transition distribution)

(Reward distribution)

The agent

The agent's policy π

$$\pi(a_t \mid s_t, \dots, s_1, a_{t-1}, \dots, a_1) \quad \text{(history-dependent policy)}$$

$$\pi(a_t \mid s_t) \quad \text{(Markov policy)}$$

Definition 4 (Utility)

Given a horizon T , the utility can be defined as

$$U_t \triangleq \sum_{k=0}^{T-t} r_{t+k} \quad (3.1)$$

The agent wants to find π **maximising** the **expected total future reward**

$$\mathbb{E}_{\mu}^{\pi} U_t = \mathbb{E}_{\mu}^{\pi} \sum_{k=0}^{T-t} r_{t+k} \quad \text{(expected utility)}$$

Contents

Bandit problems

Introduction

Bernoulli bandits

Markov processes

Markov processes

Markov decision processes and reinforcement learning

Markov decision processes

Value functions

Examples

Episodic problems

Policy evaluation

Backwards induction

Continuing, discounted problems

Markov chain theory for discounted problems

Infinite horizon MDP Algorithms

State value function

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (3.2)$$

The **optimal policy** π^*

$$\pi^*(\mu) : V_{t,\mu}^{\pi^*(\mu)}(s) \geq V_{t,\mu}^{\pi}(s) \quad \forall \pi, t, s \quad (3.3)$$

dominates all other policies π everywhere in \mathcal{S} .

The **optimal value function** V^*

$$V_{t,\mu}^*(s) \triangleq V_{t,\mu}^{\pi^*(\mu)}(s), \quad (3.4)$$

is the value function of the optimal policy π^* .

Contents

Bandit problems

Introduction

Bernoulli bandits

Markov processes

Markov processes

Markov decision processes and reinforcement learning

Markov decision processes

Value functions

Examples

Episodic problems

Policy evaluation

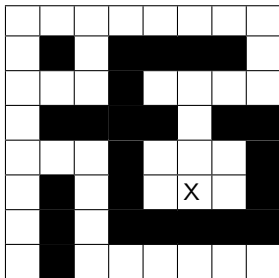
Backwards induction

Continuing, discounted problems

Markov chain theory for discounted problems

Infinite horizon MDP Algorithms

Deterministic shortest-path problems



Properties

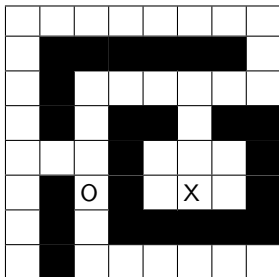
- ▶ $T \rightarrow \infty$.
- ▶ $r_t = -1$ unless $s_t = X$, in which case $r_t = 0$.
- ▶ $\mathbb{P}_\mu(s_{t+1} = X | s_t = X) = 1$.
- ▶ $\mathcal{A} = \{\text{North, South, East, West}\}$
- ▶ Transitions are deterministic and walls block.

14	13	12	11	10	9	8	7
15		13					6
16	15	14		4	3	4	5
17					2		
18	19	20		2	1	2	
19		21		1	0	1	
20		22					
21		23	24	25	26	27	28

Properties

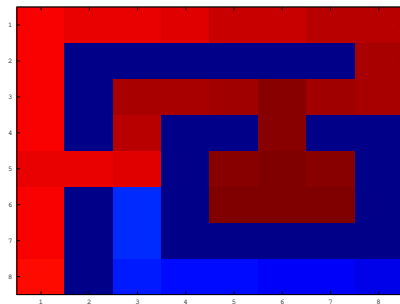
- ▶ $\gamma = 1$, $T \rightarrow \infty$.
- ▶ $r_t = -1$ unless $s_t = X$, in which case $r_t = 0$.
- ▶ The length of the shortest path from s equals the negative value of the optimal policy.
- ▶ Also called *cost-to-go*.

Stochastic shortest path problem with a pit

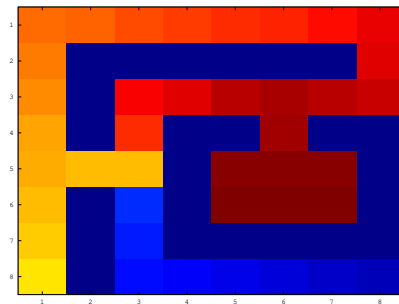


Properties

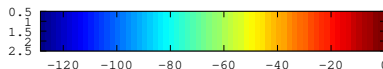
- ▶ $T \rightarrow \infty$.
- ▶ $r_t = -1$, but $r_t = 0$ at X and -100 at O and the problem ends.
- ▶ $\mathbb{P}_\mu(s_{t+1} = X | s_t = X) = 1$.
- ▶ $\mathcal{A} = \{\text{North, South, East, West}\}$
- ▶ Moves to a random direction with probability ω . Walls block.



(a) $\omega = 0.1$



(b) $\omega = 0.5$



(c) value

Figure: Pit maze solutions for two values of ω .

Exercise 3

► Why should we only take the shortcut in (a)?

► Why does the agent commit suicide at the bottom?

Contents

Bandit problems

- Introduction

- Bernoulli bandits

Markov processes

- Markov processes

Markov decision processes and reinforcement learning

- Markov decision processes

- Value functions

- Examples

Episodic problems

- Policy evaluation**

- Backwards induction

Continuing, discounted problems

- Markov chain theory for discounted problems

- Infinite horizon MDP Algorithms

How to evaluate a policy

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (4.1)$$

(4.2)

This derivation directly gives a number of **policy evaluation algorithms**.

How to evaluate a policy

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (4.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (4.2)$$

$$(4.3)$$

This derivation directly gives a number of **policy evaluation algorithms**.

How to evaluate a policy

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (4.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (4.2)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \mathbb{E}_{\mu}^{\pi}(U_{t+1} \mid s_t = s) \quad (4.3)$$

$$(4.4)$$

This derivation directly gives a number of **policy evaluation algorithms**.

How to evaluate a policy

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (4.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (4.2)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \mathbb{E}_{\mu}^{\pi}(U_{t+1} \mid s_t = s) \quad (4.3)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \sum_{i \in \mathcal{S}} V_{\mu,t+1}^{\pi}(i) \mathbb{P}_{\mu}^{\pi}(s_{t+1} = i \mid s_t = s). \quad (4.4)$$

This derivation directly gives a number of **policy evaluation algorithms**.

Monte-Carlo Policy evaluation

for $s \in \mathcal{S}$ do

end for

Monte-Carlo Policy evaluation

for $s \in \mathcal{S}$ **do**

for $k = 1, \dots, K$ **do**

 Execute policy π and record total reward K times:

$$\hat{R}_k(s) = \sum_{t=1}^T r_{t,k}.$$

end for

end for

Monte-Carlo Policy evaluation

for $s \in \mathcal{S}$ **do**

for $k = 1, \dots, K$ **do**

 Execute policy π and record total reward K times:

$$\hat{R}_k(s) = \sum_{t=1}^T r_{t,k}.$$

end for

 Calculate estimate:

$$v_t(s) = \frac{1}{K} \sum_{k=1}^K \hat{R}_k(s).$$

end for

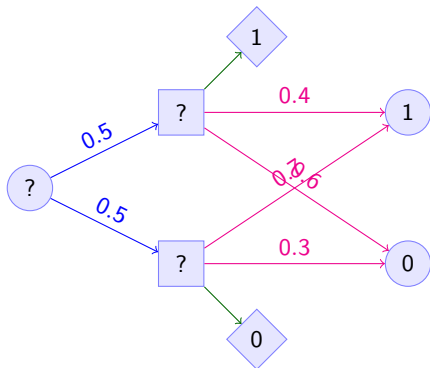
Backwards induction policy evaluation

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values of states:

$$v_t(s_t) = \sum_{a_t \in \mathcal{A}} \mathbb{P}^\pi(a_t | s_t) \left\{ \mathbb{E}_\mu(r_t | s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_\mu(s_{t+1} | s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 4

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

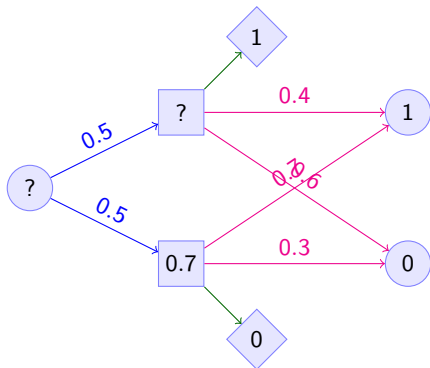
Backwards induction policy evaluation

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values of states:

$$v_t(s_t) = \sum_{a_t \in \mathcal{A}} \mathbb{P}^\pi(a_t | s_t) \left\{ \mathbb{E}_\mu(r_t | s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_\mu(s_{t+1} | s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 4

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

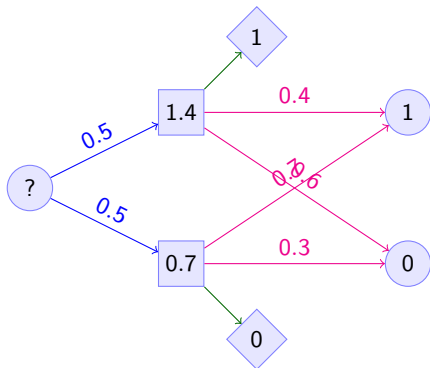
Backwards induction policy evaluation

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values of states:

$$v_t(s_t) = \sum_{a_t \in \mathcal{A}} \mathbb{P}^\pi(a_t | s_t) \left\{ \mathbb{E}_\mu(r_t | s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_\mu(s_{t+1} | s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 4

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

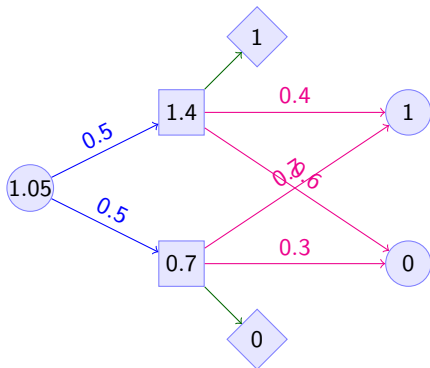
Backwards induction policy evaluation

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values of states:

$$v_t(s_t) = \sum_{a_t \in \mathcal{A}} \mathbb{P}^\pi(a_t | s_t) \left\{ \mathbb{E}_\mu(r_t | s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_\mu(s_{t+1} | s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 4

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

Contents

Bandit problems

- Introduction

- Bernoulli bandits

Markov processes

- Markov processes

Markov decision processes and reinforcement learning

- Markov decision processes

- Value functions

- Examples

Episodic problems

- Policy evaluation

- Backwards induction**

Continuing, discounted problems

- Markov chain theory for discounted problems

- Infinite horizon MDP Algorithms

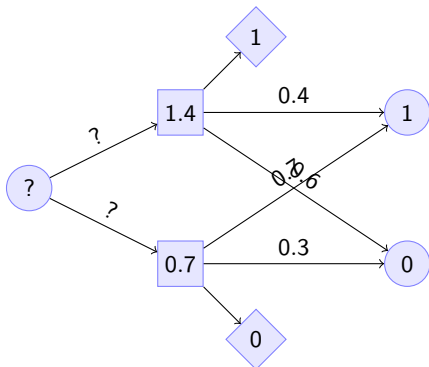
Backwards induction policy optimization

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values

$$v_t(s_t) = \max_{a_t \in \mathcal{A}} \left\{ \mathbb{E}_{\mu}(r_t \mid s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_{\mu}(s_{t+1} \mid s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 5

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

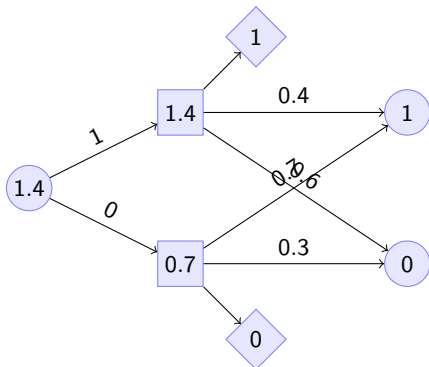
Backwards induction policy optimization

for State $s \in S$, $t = T, \dots, 1$ **do**

Update values

$$v_t(s_t) = \max_{a_t \in \mathcal{A}} \left\{ \mathbb{E}_{\mu}(r_t \mid s_t, a_t) + \sum_{s_{t+1} \in \mathcal{S}} \mathbb{P}_{\mu}(s_{t+1} \mid s_t, a_t) v_{t+1}(s_{t+1}) \right\}$$

end for



Exercise 5

What is the value $v_t(s_t)$ of the first state?

- A 1.4
- B 1.05
- C 1.0
- D 0.7
- E 0

Contents

Bandit problems

Introduction

Bernoulli bandits

Markov processes

Markov processes

Markov decision processes and reinforcement learning

Markov decision processes

Value functions

Examples

Episodic problems

Policy evaluation

Backwards induction

Continuing, discounted problems

Markov chain theory for discounted problems

Infinite horizon MDP Algorithms

Discounted total reward.

$$U_t = \lim_{T \rightarrow \infty} \sum_{k=t}^T \gamma^k r_k, \quad \gamma \in (0, 1)$$

Definition 5

A policy π is stationary if $\pi(a_t \mid s_t)$ does not depend on t .

Remark 1

We can use the Markov chain kernel $P_{\mu, \pi}$ to write the expected reward vector as

$$\mathbf{v}^\pi = \sum_{t=0}^{\infty} \gamma^t P_{\mu, \pi}^t \mathbf{r} \quad (5.1)$$

Theorem 6

For any stationary policy π , v^π is the unique solution of

$$v = r + \gamma P_{\mu, \pi} v. \quad \leftarrow \text{fixed point} \quad (5.2)$$

In addition, the solution is:

$$v^\pi = (I - \gamma P_{\mu, \pi})^{-1} r. \quad (5.3)$$

Contents

Bandit problems

- Introduction

- Bernoulli bandits

Markov processes

- Markov processes

Markov decision processes and reinforcement learning

- Markov decision processes

- Value functions

- Examples

Episodic problems

- Policy evaluation

- Backwards induction

Continuing, discounted problems

- Markov chain theory for discounted problems

- Infinite horizon MDP Algorithms

Value iteration

```
for  $n = 1, 2, \dots$  and  $s \in \mathcal{S}$  do  
     $v_n(s) = \max_a r(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_\mu(s' \mid s, a) v_{n-1}(s')$   
end for
```

Policy Iteration

```

Input  $\mu, \mathcal{S}$ .
Initialise  $v_0$ .
for  $n = 1, 2, \dots$  do
     $\pi_{n+1} = \arg \max_{\pi} \{r + \gamma P_{\pi} v_n\}$       // policy improvement
     $v_{n+1} = V_{\mu}^{\pi_{n+1}}$       // policy evaluation
    break if  $\pi_{n+1} = \pi_n$ .
end for
Return  $\pi_n, v_n$ .
  
```

- [1] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- [2] Dimitri P. Bertsekas. *Dynamic Programming and Optimal Control*. Athena Scientific, 2001.
- [3] Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- [4] Herman Chernoff. Sequential design of experiments. *Annals of Mathematical Statistics*, 30(3):755–770, 1959.
- [5] Herman Chernoff. Sequential models for clinical trials. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol.4*, pages 805–812. University of California Press, 1966.
- [6] Morris H. DeGroot. *Optimal Statistical Decisions*. John Wiley & Sons, 1970.
- [7] Marting L. Puterman. *Markov Decision Processes : Discrete Stochastic Dynamic Programming*. John Wiley & Sons, New Jersey, US, 1994.
- [8] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *ICML 2010*, 2010.
- [9] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.