

# Experiment design, Markov Decision Processes and Reinforcement Learning

## Optimal decisions, Part VII

Christos Dimitrakakis

Chalmers

November 10, 2013

## 1 Introduction

- Experiment design: examples
- Bandit problems
- Bernoulli bandits

## 2 Markov decision processes (MDP)

- Value functions

## 3 Finite horizon, undiscounted problems

- Policy evaluation
- Finite horizon backwards induction

## 4 Infinite-horizon examples

- Shortest-path problems
- Continuing problems

## 5 Infinite horizon, discounted case

- Optimality equations
- Algorithms
  - Value iteration
  - Policy iteration
  - Temporal-Difference Policy Iteration
  - Linear programming

# Clinical trials

- We have a number of treatments of unknown efficacy.
- When a new patient arrives, we must choose one of them.
- Some, slightly different, goals:
  - 1 Maximise the number of cured patients.
  - 2 Discover the best treatment.
- The optimal design is better than randomly assigning patients to treatments.

# Experimental design and Markov decision processes

The following problems

- Shortest path problems.
- Optimal stopping problems.
- Reinforcement learning problems.
- Experiment design problems.
- Multi-armed bandit problems.
- Advertising.

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

## The stochastic $n$ -armed bandit problem

- Actions  $\mathcal{A} = \{1, \dots, n\}$ .
- Expected reward  $\mathbb{E}(r_t \mid a_t = i) = \omega_i$ .
- Select actions to maximise

$$\sum_{t=0}^T \gamma^t r_t,$$

with discount factor  $\gamma \in [0, 1]$ , horizon  $T \geq 0$ .

## The stochastic $n$ -armed bandit problem

- Actions  $\mathcal{A} = \{1, \dots, n\}$ .
- Expected reward  $\mathbb{E}(r_t \mid a_t = i) = \omega_i$ .
- Select actions to maximise

$$\sum_{t=0}^T \gamma^t r_t,$$

with discount factor  $\gamma \in [0, 1]$ , horizon  $T \geq 0$ .

## Decision-theoretic approach

- Assume  $r_t \mid a_t = i \sim \psi(\omega_i)$ , with  $\omega_i \in \Omega_i$ ,  $\omega \in \Omega \triangleq \prod_i \Omega_i$  unknown parameters.
- Define prior  $\xi(\omega_1, \dots, \omega_n)$ .
- Select actions to maximise  $\mathbb{E}_\xi U_t = \mathbb{E}_\xi \sum_{k=1}^{T-t} \gamma^k r_{t+k}$ .

# Bernoulli example.

Consider  $n$  Bernoulli distributions with unknown parameters  $\omega_i$ ,  $i = 1, \dots, n$  such that

$$r_t \mid a_t = i \sim \text{Bern}(\omega_i), \quad \mathbb{E}(r_t \mid a_t = i) = \omega_i. \quad (1.1)$$

We model our belief for each bandit's parameter  $\omega_i$  as a Beta distribution  $\text{Beta}(\alpha_i, \beta_i)$ , with density  $f(\omega \mid \alpha_i, \beta_i)$  so that

$$\xi(\omega_1, \dots, \omega_n) = \prod_{i=1}^n f(\omega_i \mid \alpha_i, \beta_i).$$

$$N_{t,i} \triangleq \sum_{k=1}^t \mathbb{I}\{a_k = i\}$$

$$\hat{r}_{t,i} \triangleq \frac{1}{N_{t,i}} \sum_{k=1}^t r_t \mathbb{I}\{a_k = i\}$$

Then, the posterior distribution for the parameter of arm  $i$  is

$$\xi_t = \text{Beta}(\alpha_i + N_{t,i} \hat{r}_{t,i}, \beta_i + N_{t,i}(1 - \hat{r}_{t,i}))$$

Since  $r_t \in \{0, 1\}$  the possible states of our belief given some prior are  $\mathbb{N}^{2n}$ .

# Belief states

- The state of the bandit problem is the state of our belief.
- A sufficient statistic is the number of plays and total rewards.
- Our state  $\xi_t$  is described by the priors  $\alpha, \beta$  and the vectors

$$N_t = (N_{t,1}, \dots, N_{t,i}) \quad (1.2)$$

$$\hat{r}_t = (\hat{r}_{t,1}, \dots, \hat{r}_{t,i}). \quad (1.3)$$

- The next-state probabilities are defined as:

$$\xi_t(r_t = 1 \mid a_t = i) = \frac{\alpha_i + N_{t,i} \hat{r}_{t,i}}{\alpha_i + \beta_i + N_{t,i}}$$

- Thus decision-theoretic  $n$ -armed bandit problem can be formalised as a **Markov decision process**.



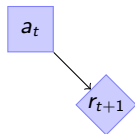


Figure: The basic bandit process

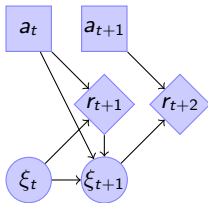


Figure: The decision-theoretic bandit process

# 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 agent obtains observations  $x_t$ .
- The agent takes actions  $a_t$  based on our observations.
- It receives rewards  $r_t$ .

## The goal (informally)

Maximise total reward  $\sum_t r_t$

## Types of environments

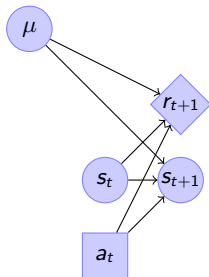
- Markov decision processes (MDPs).
- Partially observable MDPs (POMDPs).
- (Partially observable) Markov games.

# Markov decision processes

## Markov decision processes (MDP).

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

$$\mathbb{P}_{\mu}(s_{t+1} \in S \mid s_t, a_t) = \mathbb{P}_{\mu}(s_{t+1} \in S \mid s_1, a_1, \dots, s_t, a_t)$$

(Transition distribution)

$$\mathbb{P}_{\mu}(r_{t+1} \in R \mid s_t, a_t) = \mathbb{P}_{\mu}(r_{t+1} \in R \mid s_1, a_1, \dots, s_t, a_t)$$

(Reward distribution)

# The agent

## The agent's policy $\pi$

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

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

## Definition 1 (Utility)

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

We wish to find  $\pi$  maximising the expected total future reward

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

to the horizon  $T$ .

# The agent

## The agent's policy $\pi$

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

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

## Definition 1 (Utility)

$$U_t \triangleq \sum_{k=0}^{T-t} \gamma^k r_{t+k}$$

We wish to find  $\pi$  maximising the expected total future reward

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

to the horizon  $T$  with discount factor  $\gamma \in (0, 1]$ .

## State value function

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (2.1)$$

## State-action value function

$$Q_{\mu,t}^{\pi}(s, a) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s, a_t = a) \quad (2.2)$$

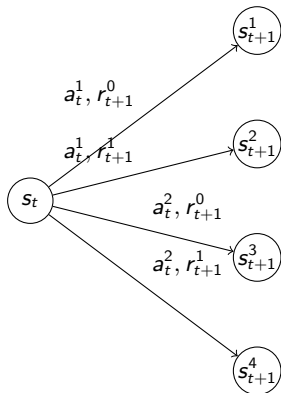
$$\pi^*(\mu) : V_{t,\mu}^{\pi^*(\mu)}(s) \geq V_{t,\mu}^{\pi}(s) \quad \forall \pi, t, s \quad (2.3)$$

The **optimal policy**  $\pi^*$  dominates all other policies  $\pi$  everywhere in  $\mathcal{S}$ .

$$V_{t,\mu}^*(s) \triangleq V_{t,\mu}^{\pi^*(\mu)}(s), \quad Q_{t,\mu}^*(s) \triangleq Q_{t,\mu}^{\pi^*(\mu)}(s, a). \quad (2.4)$$

The **optimal value function**  $V^*$  is the value function of the optimal policy  $\pi^*$ .

# Finding the optimal policy when $\mu$ is known



## Iterative/offline methods

- Estimate the optimal **value function**  $V^*$  (i.e. with backwards induction on all  $\mathcal{S}$ ).
- Iteratively **improve**  $\pi$  (i.e. with policy iteration) to obtain  $\pi^*$ .

## Online methods

- Forward **search** followed by backwards induction (on subset of  $\mathcal{S}$ ).



# Policy evaluation

## An optimal policy

*An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision. – Bellman.*

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (3.1)$$

$$(3.2)$$

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

# Policy evaluation

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (3.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (3.2)$$

$$(3.3)$$

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

# Policy evaluation

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (3.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s), \quad U_{t+1} = \sum_{k=1}^{T-t} r_{t+k}. \quad (3.2)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \mathbb{E}_{\mu}^{\pi}(U_{t+1} \mid s_t = s) \quad (3.3)$$

$$(3.4)$$

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

# Policy evaluation

$$V_{\mu,t}^{\pi}(s) \triangleq \mathbb{E}_{\mu}^{\pi}(U_t \mid s_t = s) \quad (3.1)$$

$$= \sum_{k=0}^{T-t} \mathbb{E}_{\mu}^{\pi}(r_{t+k} \mid s_t = s) \quad (3.2)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \mathbb{E}_{\mu}^{\pi}(U_{t+1} \mid s_t = s) \quad (3.3)$$

$$= \mathbb{E}_{\mu}^{\pi}(r_t \mid s_t = s) + \sum_{i \in \mathcal{S}} V_{\mu,t+1}^{\pi}(i) \mathbb{P}_{\mu}^{\pi}(s_{t+1} = i \mid s_t = s). \quad (3.4)$$

$$(3.5)$$

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

---

**Algorithm 1** Direct policy evaluation

---

```
1: for  $s \in \mathcal{S}$  do  
2:   for  $t = 0, \dots, T$  do  
3:
```

$$\hat{V}_t(s) = \sum_{k=t}^T \sum_{j \in \mathcal{S}} \mathbb{P}_{\mu}^{\pi}(s_k = j \mid s_k = s) \mathbb{E}_{\mu}^{\pi}(r_k \mid s_k = j).$$

```
4:   end for  
5: end for
```

---

**Algorithm 2** Monte-Carlo policy evaluation

---

```

for  $s \in \mathcal{S}$  do
  for  $k = 0, \dots, K$  do

```

$$\hat{V}_k(s) = \sum_{t=0}^T r_{t,k}, \quad \hat{V}(s) = \frac{1}{K} \sum_{k=1}^K \hat{V}_k(s).$$

```

  end for
end for

```

---

**Remark 1**

*The Monte Carlo evaluation algorithm has the property:*

$$\|V - \hat{V}\|_{\infty} \leq \sqrt{\frac{\ln(2|\mathcal{S}|/\delta)}{2K}}, \quad \text{with probability } 1 - \delta$$

**Proof.**

From Hoeffding's inequality, applied to any  $s$ , we have that

$$\mathbb{P} \left( |\hat{V}(s) - V(s)| \geq \sqrt{\frac{\ln(2|\mathcal{S}|/\delta)}{2K}} \right) \leq \delta/|\mathcal{S}|.$$

**Algorithm 3** Backwards induction policy evaluation

For each state  $s \in S$ , for  $t = 1, \dots, T - 1$ :

$$\hat{V}_t(s) = r(s) + \sum_{j \in S} \mathbb{P}_{\mu, \pi}(s_{t+1} = j \mid s_t = s) \hat{V}_{t+1}(j), \quad (3.6)$$

with  $\hat{V}_T(s) = r(s)$ .

**Theorem 2**

*Algorithm 3 results in estimates with the property:*

$$\hat{V}_t(s) = V_{\mu, t}^{\pi}(s) \quad (3.7)$$

**Algorithm 4** Finite-horizon backwards induction

Input  $\mu, \mathcal{S}_T$ .

Initialise  $V_T(s)$ , for all  $s \in \mathcal{S}_T$ .

**for**  $n = T - 1, T - 2, \dots, 1$  **do**

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

$$\pi_n(s) = \arg \max_a \mathbb{P}_\mu(s'|s, a) [\mathbb{E}_\mu(r|s', s) + V_{n+1}^*(s')]$$

$$V_n(s) = \sum_{s' \in \mathcal{S}_{n+1}} \mathbb{P}_\mu(s'|s, \pi_n(s)) [\mathbb{E}_\mu(r|s', s) + V_{n+1}(s')]$$

**end for**

**end for**

Return  $\pi = (\pi_n)_{n=1}^T$ .

**Notes**

- $\mathbb{P}_{\mu, \pi}(s'|s) = \sum_a \mathbb{P}_\mu(s'|s, a) \mathbb{P}_\pi(a|s)$ .
- Finite horizon problems only, or approximations (e.g. lookahead in game trees).
- For stochastic problems, we marginalize over states.
- As we know the optimal choice at the last step, we can find the optimal policy!
- Can be used with estimates of the value function.



### Theorem 3

*For a  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

### Proof.

**1** First we show that  $V_t \geq V_t^*$ .



## Theorem 3

*For a  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .



## Theorem 3

*For a  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .
- 3 Assume that for  $n \geq t + 1$ , (3.8) holds.



## Theorem 3

*For a  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .
- 3 Assume that for  $n \geq t + 1$ , (3.8) holds.
- 4 Then it holds for  $n = t$  since:

$$V_t(s) = \max_a \left\{ r(s) + \sum_{j \in \mathcal{S}} p(j|s, a) V_{t+1}(j) \right\}$$

## Theorem 3

For a  $T$ -horizon problems, backwards induction is optimal, i.e.

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .
- 3 Assume that for  $n \geq t + 1$ , (3.8) holds.
- 4 Then it holds for  $n = t$  since:

$$V_t(s) \geq \max_a \left\{ r(s) + \sum_{j \in \mathcal{S}} p(j|s, a) V_{\mu,t+1}^*(j) \right\} \quad (\text{by step 3})$$

## Theorem 3

*For a  $T$ -horizon problems, backwards induction is optimal, i.e.*

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .
- 3 Assume that for  $n \geq t+1$ , (3.8) holds.
- 4 Then it holds for  $n = t$  since:

$$V_t(s) \geq \max_a \left\{ r(s) + \sum_{j \in \mathcal{S}} p(j|s, a) V_{\mu,t+1}^{\pi'}(j) \right\}$$

## Theorem 3

For a  $T$ -horizon problems, backwards induction is optimal, i.e.

$$V_n(s) = V_{\mu,n}^*(s) \quad (3.8)$$

## Proof.

- 1 First we show that  $V_t \geq V_t^*$ .
- 2 For  $n = T$ ,  $V_T(s) = r(s) = V_{\mu,T}^\pi(s)$ .
- 3 Assume that for  $n \geq t + 1$ , (3.8) holds.
- 4 Then it holds for  $n = t$  since:

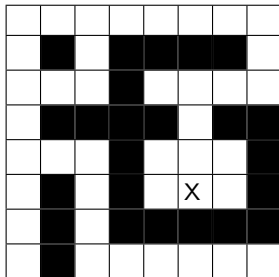
$$V_t(s) \geq V_t^{\pi'}(s)$$

- 5 The above holds for any policy  $\pi'$ , including  $\pi' = \pi$ , the policy returned by backwards induction. Then:

$$V_{\mu,t}^*(s) \geq V_{\mu,t}^\pi(s) = V_t(s) \geq V_{\mu,t}^*(s).$$



# Deterministic shortest-path problems



## Properties

- $\gamma = 1, T \rightarrow \infty$ .
- $r_t = -1$  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.

What is the shortest path to the destination from any point?



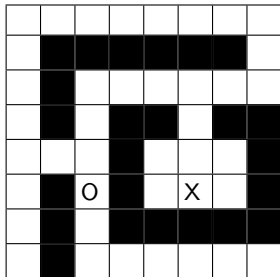
# Shortest-path problem solution

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*.
- Remember Dijkstra's algorithm?

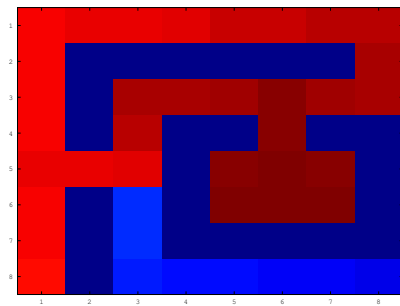
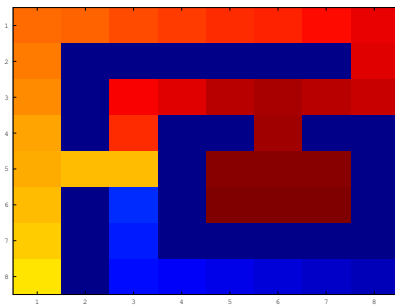
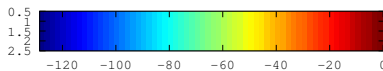
## Stochastic shortest path problem, with a pit



## Properties

- $\gamma = 1$ ,  $T \rightarrow \infty$ .
- $r_t = -1$ , but  $r_t = 0$  at X and  $-100$  at O and episode 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  $\omega$ . Walls block.

For what value of  $\omega$  is it better to take the dangerous shortcut? (However, if we want to take into account risk explicitly we must modify the agent's utility function)

(a)  $\omega = 0.1$ (b)  $\omega = 0.5$ 

(c) value

Figure: Pit maze solutions for two values of  $\omega$ .

# Continuing stochastic MDPs

## Inventory management

- There are  $K$  storage locations.
- Each place can store  $n_i$  items.
- At each time-step there is a probability  $\phi_i$  that a client try to buy an item from location  $i$ ,  $\sum_i \phi_i \leq 1$ . If there is an item available, you gain reward 1.
- Action 1: ordering  $u$  units of stock, for paying  $c(u)$ .
- Action 2: move  $u$  units of stock from one location  $i$  to another,  $j$ , for a cost  $\psi_{ij}(u)$ .

## An easy special case

- $K = 1$ .
- There is one type of item only.
- Orders are placed and received every  $n$  timesteps.

# Inventory management

## An easy special case

- $K = 1$ .
- Deliveries happen once every  $m$  timesteps.
- Each time-step a client arrives with probability  $\phi$ .

## Properties

- The state set .
- The action set .
- The transition probabilities

# Inventory management

## An easy special case

- $K = 1$ .
- Deliveries happen once every  $m$  timesteps.
- Each time-step a client arrives with probability  $\phi$ .

## Properties

- The state set is the number of items we have:  $\mathcal{S} = \{0, 1, \dots, n\}$ .
- The action set .
- The transition probabilities

# Inventory management

## An easy special case

- $K = 1$ .
- Deliveries happen once every  $m$  timesteps.
- Each time-step a client arrives with probability  $\phi$ .

## Properties

- The state set is the number of items we have:  $\mathcal{S} = \{0, 1, \dots, n\}$ .
- The action set  $\mathcal{A} = \{0, 1, \dots, n\}$  since we can order from nothing up to  $n$  items.
- The transition probabilities

# Inventory management

## An easy special case

- $K = 1$ .
- Deliveries happen once every  $m$  timesteps.
- Each time-step a client arrives with probability  $\phi$ .

## Properties

- The state set is the number of items we have:  $\mathcal{S} = \{0, 1, \dots, n\}$ .
- The action set  $\mathcal{A} = \{0, 1, \dots, n\}$  since we can order from nothing up to  $n$  items.
- The transition probabilities  $P(s'|s, a) = \binom{m}{d} \phi^d (1 - \phi)^{m-d}$ , where  $d = s + a - s'$ , for  $s + a \leq n$ .



## Discounted total reward

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

## Definition 4

A policy  $\pi$  is stationary if  $\pi(a_t \mid s_t) = \pi(a_n \mid s_n)$  for all  $n, t$ .

## Remark 2

We can use the Markov chain kernel  $\mathbf{P}$  to write the expected reward vector as

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

## Theorem 5

For any stationary  $\pi$ ,  $\mathbf{v}^\pi$  is the unique solution of

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

In addition, the solution is:

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

Proof.



## Theorem 5

For any stationary  $\pi$ ,  $\mathbf{v}^\pi$  is the unique solution of

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

In addition, the solution is:

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

Proof.

$$\mathbf{1} \quad \mathbf{r} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \mathbf{v}$$



## Theorem 5

For any stationary  $\pi$ ,  $\mathbf{v}^\pi$  is the unique solution of

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

In addition, the solution is:

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

## Proof.

1  $\mathbf{r} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \mathbf{v}$

2 Since  $\|\gamma \mathbf{P}_{\mu, \pi}\| < 1 \cdot \|\mathbf{P}_{\mu, \pi}\| = 1$ , the following inverse exists:

$$(\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} = \lim_{n \rightarrow \infty} \sum_{t=0}^n (\gamma \mathbf{P}_{\mu, \pi})^t$$

## Theorem 5

For any stationary  $\pi$ ,  $\mathbf{v}^\pi$  is the unique solution of

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

In addition, the solution is:

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

## Proof.

**1**  $\mathbf{r} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi}) \mathbf{v}$

**2** Since  $\|\gamma \mathbf{P}_{\mu, \pi}\| < 1 \cdot \|\mathbf{P}_{\mu, \pi}\| = 1$ , the following inverse exists:

$$(\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} = \lim_{n \rightarrow \infty} \sum_{t=0}^n (\gamma \mathbf{P}_{\mu, \pi})^t$$

**3** Using step 1 and then 2,

$$\mathbf{v} = (\mathbf{I} - \gamma \mathbf{P}_{\mu, \pi})^{-1} \mathbf{r} = \sum_{t=0}^{\infty} \gamma^t \mathbf{P}_{\mu, \pi}^t \mathbf{r} = \mathbf{v}^\pi,$$

where the last step is by Remark 2



## Definition 6 (Bellman operator)

$$\mathcal{L}_\pi \mathbf{v} \triangleq \mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v}$$

$$\mathcal{L} \mathbf{v} \triangleq \sup_{\pi} \{ \mathbf{r} + \gamma \mathbf{P}_\pi \mathbf{v} \}, \quad \mathbf{v} \in \mathcal{V}$$

$$\mathbf{v} = \mathcal{L} \mathbf{v} \quad (\text{Bellman optimality equation})$$

## Theorem 7

For any bounded  $\mathbf{r}$ , it holds that for  $\mathbf{v} \in \mathcal{V}$ :

- If  $\mathbf{v} \geq \mathcal{L} \mathbf{v}$ , then  $\mathbf{v} \geq \mathbf{v}^*$
- If  $\mathbf{v} \leq \mathcal{L} \mathbf{v}$ , then  $\mathbf{v} \leq \mathbf{v}^*$
- If  $\mathbf{v} = \mathcal{L} \mathbf{v}$ , then  $\mathbf{v}$  is unique and  $\mathbf{v} = \mathbf{v}^*$ ,

where  $\mathbf{v}^* = \sup_{\pi} \mathbf{v}^{\pi}$ .

## Theorem 8 (Banach Fixed-Point theorem)

Suppose  $S$  is a Banach space (i.e. a complete normed linear space) and  $T : S \rightarrow S$  is a contraction mapping (i.e.  $\exists \gamma \in [0, 1)$  s.t.  $\|Tu - Tv\| \leq \gamma\|u - v\|$  for all  $u, v \in S$ ). Then

- There is a unique  $u^* \in U$  s.t.  $Tu^* = u^*$  and
- For any  $u^0 \in S$  the sequence  $\{u^n\}$ :

$$u^{n+1} = Tu^n = T^{n+1}u^0$$

converges to  $u^*$ .

## Proof.

For any  $m \geq 1$

$$\|u^{n+m} - u^n\| \leq \sum_{k=0}^{m-1} \|u^{n+k+1} - u^{n+k}\| = \sum_{k=0}^{m-1} \|T^{n+k}u^1 - T^{n+k}u^0\|$$



## Theorem 8 (Banach Fixed-Point theorem)

Suppose  $S$  is a Banach space (i.e. a complete normed linear space) and  $T : S \rightarrow S$  is a contraction mapping (i.e.  $\exists \gamma \in [0, 1)$  s.t.  $\|Tu - Tv\| \leq \gamma \|u - v\|$  for all  $u, v \in S$ ). Then

- There is a unique  $u^* \in S$  s.t.  $Tu^* = u^*$  and
- For any  $u^0 \in S$  the sequence  $\{u^n\}$ :

$$u^{n+1} = Tu^n = T^{n+1}u^0$$

converges to  $u^*$ .

## Proof.

For any  $m \geq 1$

$$\begin{aligned} \|u^{n+m} - u^n\| &\leq \sum_{k=0}^{m-1} \|u^{n+k+1} - u^{n+k}\| = \sum_{k=0}^{m-1} \|T^{n+k}u^1 - T^{n+k}u^0\| \\ &\leq \sum_{k=0}^{m-1} \gamma^{n+k} \|u^1 - u^0\| = \frac{\gamma^n(1 - \gamma^m)}{1 - \gamma} \|u^1 - u^0\|. \end{aligned}$$





## Theorem 9

If  $\gamma \in [0, 1)$  then the Bellman operator  $\mathcal{L}$  is a contraction mapping in  $\mathcal{V}$ .

### Proof.

Let  $\mathbf{v}, \mathbf{v}' \in \mathcal{V}$ . Consider  $s \in \mathcal{S}$  s. t.  $\mathcal{L}\mathbf{v}(s) \geq \mathcal{L}\mathbf{v}'(s)$ , and let

$$a_s^* \in \arg \max_{a \in \mathcal{A}} \left\{ r(s) + \sum_{j \in \mathcal{S}} \gamma p_{\mu}(j \mid s, a) \mathbf{v}(j) \right\}.$$

Then

$$0 \leq \mathcal{L}\mathbf{v}(s) - \mathcal{L}\mathbf{v}'(s) \leq \gamma \|\mathbf{v} - \mathbf{v}'\|.$$

Repeating the argument for  $s$  such that  $\mathcal{L}\mathbf{v}(s) \leq \mathcal{L}\mathbf{v}'(s)$ , we obtain

$$|\mathcal{L}\mathbf{r}(s) - \mathcal{L}\mathbf{r}'(s)| \leq \gamma \|\mathbf{r} - \mathbf{r}'\|.$$

Taking the supremum, we obtain the required result. □

## Theorem 10

If  $\gamma \in [0, 1)$ ,  $S$  is discrete and  $r$  is bounded:

- There is a unique  $v^* \in \mathcal{V}$  s.t.  $\mathcal{L}v^* = v^*$  and such that  $v^* = V_\mu^*$ .
- For a stationary  $\pi$ , there is a unique  $v \in \mathcal{V}$  such that  $\mathcal{L}_\pi v = v$  and  $v = V_\mu^\pi$ .

## Proof.

- From the previous theorem,  $\mathcal{L}$  is a contraction. So, we can apply the Fixed-Point theorem. Thus there is a unique solution. This is the optimal value function due to Theorem 9
- Use part 1 with  $\Pi = \{\pi\}$ .



---

**Algorithm 5** Value iteration

---

Input  $\mu, \mathcal{S}$ .

Initialise  $v_0 \in \mathcal{V}$ .

**for**  $n = 1, 2, \dots$  **do**

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

$$\pi_n(s) = \arg \max_a r(s) \sum_{s' \in \mathcal{S}} \mathbb{P}_\mu(s'|s, a) v_{n-1}(s')$$

$$v_n(s) = r(s) + \sum_{s' \in \mathcal{S}} \mathbb{P}_\mu(s'|s, \pi_n(s)) v_{n-1}(s')$$

**end for**

**break** if termination-condition is met

**end for**

Return  $\pi_n, V_n$ .

---

## Theorem 11

The value iteration algorithm satisfies

- $\lim_{n \rightarrow \infty} \|v_n - V^*\| = 0.$
- There exists  $N < \infty$  such that

$$\|v_{n+1} - v_n\| \leq \epsilon(1 - \gamma)/2\gamma, \quad \forall n \geq N. \quad (5.4)$$

- The policy  $\pi_\epsilon$  that takes action  $\arg \max_a r(s) + \gamma \sum_j p(j|s, a)v_n(s')$  is  $\epsilon$ -optimal.
- $\|v_{n+1} - V_\mu^*\| < \epsilon/2$  for  $n > N$ .

## Proof.

The first two statements follow from the fixed point theorem. Now note that

$$\|V_\mu^{\pi_\epsilon} - V_\mu^*\| \leq \|V_\mu^{\pi_\epsilon} - v_n\| + \|v_n - V_\mu^*\|$$

We can bound these two terms easily:

$$\|V^{\pi_\epsilon} - v_{n+1}\| \leq \frac{\gamma}{1 - \gamma} \|v_{n+1} - v_n\|, \quad \|v_{n+1} - V_\mu^*\| \leq \frac{\gamma}{1 - \gamma} \|v_{n+1} - v_n\|$$



## Theorem 12

*Value iteration converges linearly at rate  $\gamma$  and  $O(\gamma^n)$ . In addition, for  $r \in [0, 1]$  and  $r^0 = \mathbf{0}$*

$$\|v_n - V_\mu^*\| \leq \frac{\gamma^n}{1 - \gamma}$$
$$\|V_\mu^{\pi_n} - V_\mu^*\| \leq \frac{2\gamma^n}{1 - \gamma},$$

---

**Algorithm 6** Policy iteration

---

Input  $\mu, \mathcal{S}$ .

Initialise  $v_0$ .

**for**  $n = 1, 2, \dots$  **do**

$\pi_{n+1} = \arg \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_n$  (policy improvement)

$\mathbf{v}_{n+1} = V_{\mu}^{\pi_{n+1}}$  (policy evaluation)

**break** if  $\pi_{n+1} = \pi_n$ .

**end for**

Return  $\pi_n, \mathbf{v}_n$ .

---

## Theorem 13

*If  $v_n, v_{n+1}$  are produced by policy iteration, then  $v_n \leq v_{n+1}$ .*

### Proof.

From the policy improvement step

$$r + \gamma P_{\pi_{n+1}} v_n \geq r + \gamma P_{\pi_n} v_n = v_n$$

where the equality is due to the fact that  $(I - \gamma P_{\mu, \pi_n})v_n = r$  from the policy evaluation step. Rearranging, we get that

$$\begin{aligned} r &\geq (I - \gamma P_{\pi_{n+1}})v_n \\ (I - \gamma P_{\pi_{n+1}})^{-1}r &\geq v_n, \end{aligned}$$

noting that the inverse is positive. Since the left side equals  $v_{n+1}$ , we have proved the theorem. □

## Corollary 14

*If  $S, \mathcal{A}$  are finite then policy iteration terminates in a finite number of iterations.*

# Modified policy iteration

---

## Algorithm 7 Modified policy iteration

---

Input  $\mu, \mathcal{S}$ .

Initialise  $v_0$ .

**for**  $n = 1, 2, \dots$  **do**

$\pi_n = \arg \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_{n-1}$  // policy improvement

$\mathbf{v}_n = \mathcal{L}_{\pi_n}^k \mathbf{v}_{n-1}$  // partial policy evaluation

**break** if  $\pi_n = \pi_{n-1}$ .

**end for**

Return  $\pi_n, \mathbf{v}_n$ .

---



# Geometric view

## Definition 15

Difference operator

$$\mathcal{B}v \triangleq \max_{\pi} \{r + (\gamma P_{\pi} - I)v\} = \mathcal{L}v - v. \quad (5.5)$$

Hence the optimality equation becomes

$$\mathcal{B}v = 0. \quad (5.6)$$

.

---

**Algorithm 8** Temporal-Difference Policy Iteration
 

---

Input  $\mu, \mathcal{S}, \lambda$ .

Initialise  $\mathbf{v}_0$ .

**for**  $n = 1, 2, \dots$  **do**

$\pi_n = \arg \max_{\pi} \mathbf{r} + \gamma \mathbf{P}_{\pi} \mathbf{v}_{n-1}$  // policy improvement

$\mathbf{v}_n = \mathbf{v}_{n-1} + \tau_k$  // temporal difference update.

**break** if  $\pi_n = \pi_{n=1}$ .

**end for**

Return  $\pi_n, \mathbf{v}_n$ .

---

$$\mathcal{L}_{\pi_{n+1}} \mathbf{v}_n = \mathcal{L} \mathbf{v}_n. \quad (5.7)$$

$$d_n(i, j) = \mathbf{v}_n(i) - [\mathbf{r}(i) + \gamma \mathbf{v}_n(j)]. \quad (\text{temporal difference error})$$

$$\tau_n(i) = \sum_{t=0}^{\infty} \mathbb{E}_{\pi_n, \mu} [(\gamma \lambda)^t d_n(s_t, s_{t+1}) \mid s_0 = i] \quad (5.8)$$

$$\mathbf{v}_{n+1} = \mathbf{v}_n + \tau_n \quad (5.9)$$

$$\mathcal{D}_n \mathbf{v} \triangleq (1 - \lambda) \mathcal{L}_{\pi_{n+1}} \mathbf{v}_n + \lambda \mathcal{L}_{\pi_{n+1}} \mathbf{v}, \quad (\text{fixed point})$$

Select  $\mathbf{y} \in \mathbb{S}^{|\mathcal{S}|}$  (i.e. a state distribution). Then:

### Primal linear program

$$\min_{\mathbf{v}} \mathbf{y}^\top \mathbf{v}$$

such that

$$\mathbf{v}(s) - \gamma \mathbf{p}_{s,a}^\top \mathbf{v} \geq r(s, a), \quad \forall a \in \mathcal{A}, s \in \mathcal{S}.$$

### Dual linear program

$$\max_x \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} r(s, a) x(s, a)$$

such that  $x \in \mathbb{R}_+^{|\mathcal{S} \times \mathcal{A}|}$  and

$$\sum_{a \in \mathcal{A}} x(j, a) - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \gamma p(j|s, a) x(s, a) = y(j).$$

with  $\mathbf{y} \in \mathbb{S}^{|\mathcal{S}|}$ .

# Summary

## Markov decision processes

Can represent : Shortest path problems, Stopping problems, Experiment design problems, Multi-armed bandit problems, Reinforcement learning problems.

## Backwards induction (aka value iteration)

- In the class of dynamic programming algorithms.
- Tractable when either the state space  $\mathcal{S}$  or the horizon  $T$  are small (finite).

## Optimal decisions and Bayesian reinforcement learning

- A known environment is represented as an MDP.
- Bandit problems can be solved by representing them as infinite-state MDPs.
- In general, an unknown environment can be represented as a distribution over MDPs.
- The decision problem can again be formulated as an infinite-state MDP.

- [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. Univ. of Calif 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.