

Markov Decision Processes and Policy Iteration

Last Time

- What does "**Markov**" mean in "Markov Process"?
- What is a **Markov decision process**?

Guiding Questions

Guiding Questions

- What is a **Markov decision process**?

Guiding Questions

- What is a **Markov decision process**?
- What is a **policy**?

Guiding Questions

- What is a **Markov decision process**?
- What is a **policy**?
- How do we **evaluate** policies?

MDP "Tuple Definition"

MDP "Tuple Definition"

$$(S, A, T, R, \gamma)$$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states $\{1, 2, 3\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
- $\{1, 2, 3\}$
 $\{\text{healthy, pre-cancer, cancer}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states $\{1, 2, 3\}$ \mathbb{R}^2
 $\{\text{healthy, pre-cancer, cancer}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
- $\{1, 2, 3\}$ \mathbb{R}^2 $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
- $\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
- $\{1, 2, 3\} \quad (x, y) \in \mathbb{R}^2 \quad \{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\} \quad (s, i, r) \in \mathbb{N}^3$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
 $\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$ $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
 $\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$ $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions $\{1, 2, 3\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
 $\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$ $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions
 $\{1, 2, 3\}$
 $\{\text{test, wait, treat}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states
 - $\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 - $\{\text{healthy, pre-cancer, cancer}\}$ $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions
 - $\{1, 2, 3\}$ \mathbb{R}^2
 - $\{\text{test, wait, treat}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$

$\{\text{test, wait, treat}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- T (transition distribution) - explicit or implicit ("generative")
model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$

$\{\text{test, wait, treat}\}$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$

$\{\text{test, wait, treat}\}$
- T (transition distribution) - explicit or implicit ("generative")
model of how the state changes

$T(s' \mid s, a)$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$
- R (reward function) - maps each state and action to a reward

$\{\text{test, wait, treat}\}$
 $T(s' \mid s, a)$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$
- R (reward function) - maps each state and action to a reward

$\{\text{test, wait, treat}\}$
 $T(s' \mid s, a)$

$R(s, a)$ or
 $R(s, a, s')$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
 - A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
 - T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$
 - R (reward function) - maps each state and action to a reward

$\{\text{test, wait, treat}\}$
 $T(s' \mid s, a)$
- $R(s, a)$ or
 $R(s, a, s')$
- $s', r = G(s, a)$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$
 $\{\text{test, wait, treat}\}$
 $T(s' \mid s, a)$
- R (reward function) - maps each state and action to a reward

$R(s, a)$ or
 $R(s, a, s')$
- γ : discount factor

$s', r = G(s, a)$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b and/or S_T in some contexts)

- S (state space) - set of all possible states

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- A (action space) - set of all possible actions

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$\{1, 2, 3\}$
 \mathbb{R}^2
 $\{0, 1\} \times \mathbb{R}^2$

 $T(s' \mid s, a)$
- R (reward function) - maps each state and action to a reward

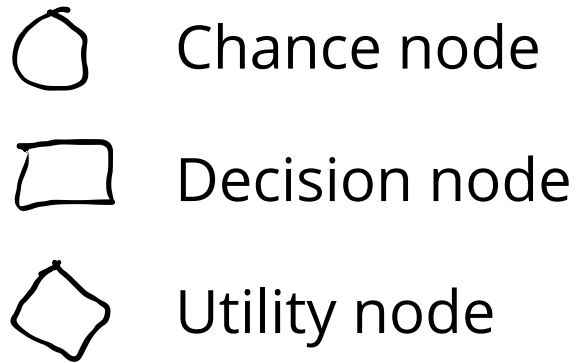
$\{\text{test, wait, treat}\}$

 $R(s, a)$ or $R(s, a, s')$
- γ : discount factor
- b : initial state distribution
- S_t : set of terminal states

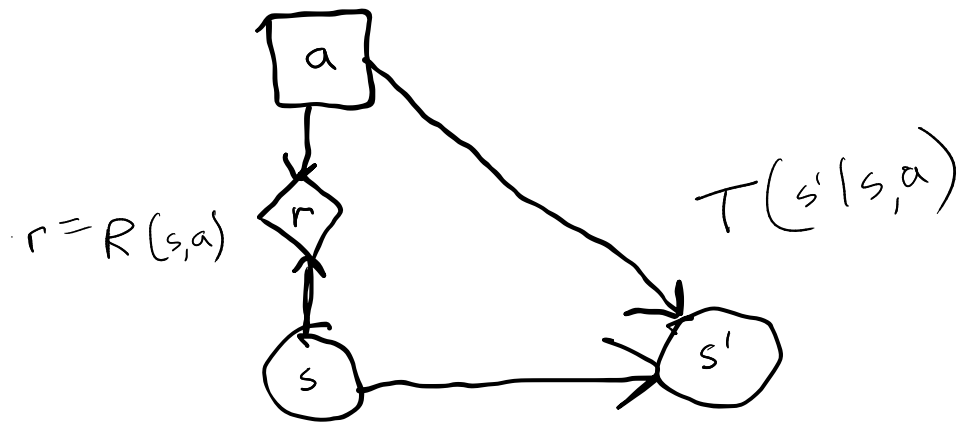
$$s', r = G(s, a)$$

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network



S = set all values that s and s' can take
 A = " " " " a can take
 T
 R
 γ

MDP Example

Imagine it's a cold day and you're ready to go to work. You have to decide whether to bike or drive.

MDP Example

Imagine it's a cold day and you're ready to go to work. You have to decide whether to bike or drive.

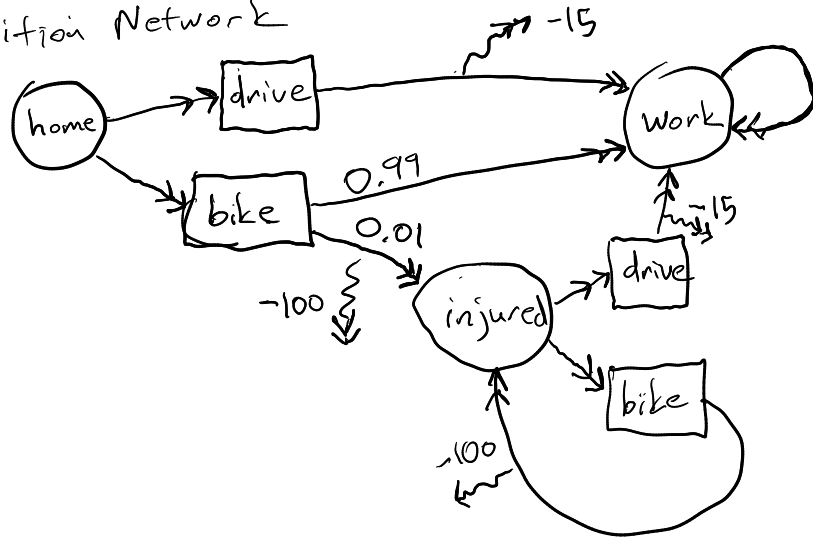
- If you drive, you will have to pay \$15 for parking; biking is free.

MDP Example

Imagine it's a cold day and you're ready to go to work. You have to decide whether to bike or drive.

- If you drive, you will have to pay \$15 for parking; biking is free.
- On 1% of cold days, the ground is covered in ice and you will crash if you bike, but you can't discover this until you start riding. After your crash, you limp home with pain equivalent to losing \$100.

Not a Decision Network
Transition Network



$$S = \{\text{home}, \text{injured}, \text{work}\}$$

$$A = \{\text{drive}, \text{bike}\}$$

$$T^{\text{drive}} = \begin{matrix} & \begin{matrix} h & i & w \end{matrix} \\ \begin{matrix} h \\ i \\ w \end{matrix} & \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \end{matrix} \quad T^{\text{bike}} = \begin{bmatrix} 0 & 0.01 & 0.99 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R(s, a, s') = \begin{cases} -15 & \text{if } a = \text{drive} \\ -100 & \text{if } s' = \text{injured} \\ 0 & \text{o.w.} \end{cases}$$

$$\gamma = 0.99$$

$$S_T = \{\text{work}\} \quad b(s) = \begin{cases} 1 & \text{if } s = \text{home} \\ 0 & \text{o.w.} \end{cases}$$

Policies and Simulation

Policies and Simulation

- A *policy*, denoted with $\pi(a_t \mid s_t)$, is a conditional distribution of actions given states.

Policies and Simulation

- A *policy*, denoted with $\pi(a_t \mid s_t)$, is a conditional distribution of actions given states.
- $a_t = \pi(s_t)$ is used as shorthand when a policy is deterministic.

Policies and Simulation

- A *policy*, denoted with $\pi(a_t \mid s_t)$, is a conditional distribution of actions given states.
- $a_t = \pi(s_t)$ is used as shorthand when a policy is deterministic.
- When a policy is combined with an MDP, it becomes a Markov stochastic process with

$$P(s' \mid s) = \sum_a T(s' \mid s, a) \pi(a \mid s)$$

Policies and Simulation

- A *policy*, denoted with $\pi(a_t \mid s_t)$, is a conditional distribution of actions given states.
- $a_t = \pi(s_t)$ is used as shorthand when a policy is deterministic.
- When a policy is combined with an MDP, it becomes a Markov stochastic process with

$$P(s' \mid s) = \sum_a T(s' \mid s, a) \pi(a \mid s)$$

MDP Objective:

$$U(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi \right]$$

Policies and Simulation

- A *policy*, denoted with $\pi(a_t \mid s_t)$, is a conditional distribution of actions given states.
- $a_t = \pi(s_t)$ is used as shorthand when a policy is deterministic.
- When a policy is combined with an MDP, it becomes a Markov stochastic process with

$$P(s' \mid s) = \sum_a T(s' \mid s, a) \pi(a \mid s)$$

MDP Objective:

$$\underline{U(\pi)} = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi \right]$$

Algorithm: Rollout Simulation

Inputs: MDP (S, A, R, T, γ, b) (only need generative model, G), Policy π , horizon H

Outputs: Utility estimate \hat{u}

$s \leftarrow \text{sample}(b)$

$\hat{u} \leftarrow 0$

for t in $0 \dots H - 1$

$a \leftarrow \text{sample}(\pi(a \mid s))$

$s', r \leftarrow G(s, a)$.

$\hat{u} \leftarrow \hat{u} + \gamma^t r$

$s \leftarrow s'$

return \hat{u}

Policy Evaluation

Naive Policy Evaluation not on Exam

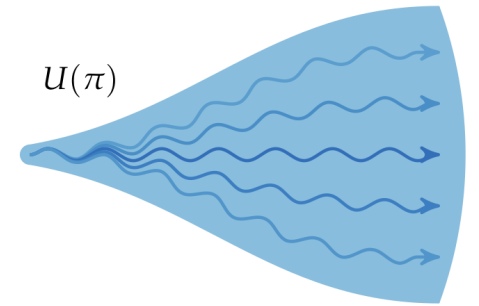
Monte Carlo Policy Evaluation

Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

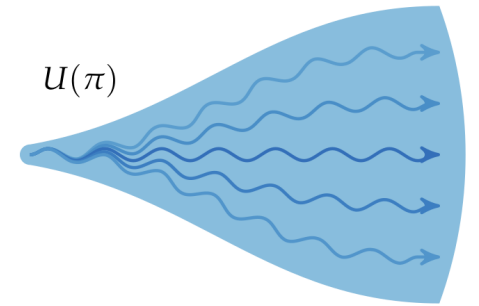
Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



Monte Carlo Policy Evaluation

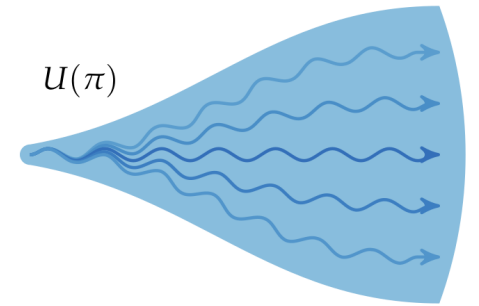
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

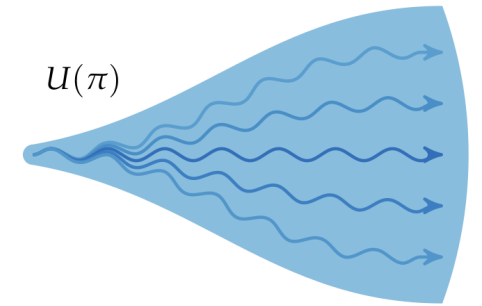


Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



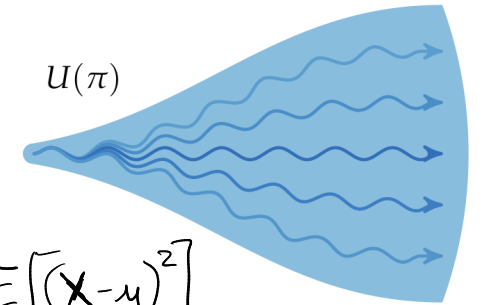
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$
$$U(\pi) \approx \bar{u}_m \overset{\substack{\text{R.V. itself} \\ \swarrow}}{=} \frac{1}{m} \sum_{i=1}^m \hat{u}^{(i)}$$

where $\hat{u}^{(i)}$ is generated by a rollout simulation

Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

$$\sigma_x = \text{standard dev} \quad \mu = E[X]$$

$$\sigma_x^2 = \text{Var}(X) = E[(X - \mu)^2]$$

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

$$U(\pi) \approx \bar{u}_m = \frac{1}{m} \sum_{i=1}^m \hat{u}^{(i)} \quad \text{Var}(\bar{u}_m) = \frac{\sigma^2}{m} \quad \text{Var}(\hat{u}) = \sigma^2$$

How can we quantify the accuracy of \bar{u}_m ?

$$\begin{aligned} \text{Var}(\bar{u}_m) &= \text{Var}\left(\frac{1}{m} \sum_i \hat{u}^{(i)}\right) \\ &= \frac{1}{m^2} \text{Var}\left(\sum_i \hat{u}^{(i)}\right) \\ &= \frac{1}{m^2} \sum_{i=1}^m \text{Var}(\hat{u}^{(i)}) \quad (\text{Bienaymé}) \end{aligned}$$

$$\frac{\sigma^2}{m} = \frac{1}{m^2} m \hat{\sigma}^2 \Rightarrow \bar{\sigma} = \frac{\hat{\sigma}}{\sqrt{m}}$$

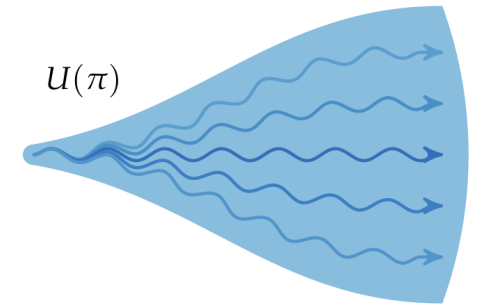
where $\hat{u}^{(i)}$ is generated by a rollout simulation

Standard Error of Mean

$$S.E.M. \equiv \frac{\text{std}(\hat{u})}{\sqrt{m}}$$

Monte Carlo Policy Evaluation

- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$

How can we quantify the accuracy of \bar{u}_m ?

also an R.V.

$$U(\pi) \approx \bar{u}_m = \frac{1}{m} \sum_{i=1}^m \hat{u}^{(i)}$$

where $\hat{u}^{(i)}$ is generated by a rollout simulation

Value Function-Based Policy Evaluation

Discrete, finite, state and action spaces

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid \pi, s_0 = s \right]$$

$$U(\pi) = E_{s \sim b} [U^\pi(s)]$$

$$= E[r_0 \mid s_0 = s, \pi] + E \left[\sum_{t=1}^{\infty} \gamma^t r_t \mid \pi, s_0 = s \right]$$

$$= R(s, \pi(s)) + \sum_{s' \in S} T(s' \mid s, a) E \left[\sum_{t=1}^{\infty} \gamma^t r_t \mid \pi, s_0 = s, s_1 = s' \right]$$

$\tau = t-1$

$$= R(s, \pi(s)) + \gamma \sum_{s'} T(s' \mid s, a) E \left[\sum_{\tau=0}^{\infty} \gamma^\tau r_\tau \mid \pi, s_0 = s' \right]$$

$U^\pi(s')$

$$U^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s' \mid s, a) U^\pi(s')$$

Bellman's Expectation Eq.

$$\vec{U}^\pi \quad \vec{U}_{\text{ind}(s)}^\pi = U^\pi(s)$$

$$\vec{R}_{\text{ind}(s)}^\pi = R(s, \pi(s))$$

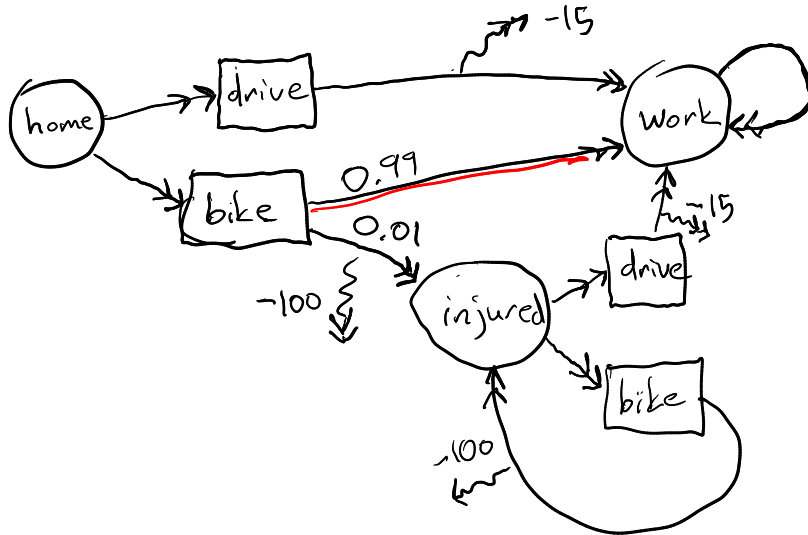
$$T_{\text{ind}(s), \text{ind}(s')}^\pi = T(s' \mid s, a)$$

$$\vec{U}^\pi = \vec{R}^\pi + \gamma T^\pi \vec{U}^\pi$$

$$\vec{U}^\pi = (\mathbf{I} - \gamma T^\pi)^{-1} \vec{R}^\pi$$

Break

- Suggest a policy that you think is optimal for the icy day problem



$$U(\text{drive from home}) = -15$$

$$U(\text{bike from home}) = 0.01 \times (-100 + \overset{0.9}{\downarrow} -15) = -1.15$$

$$\pi(s) = \begin{cases} \text{bike} & \text{if } s = \text{home} \\ \text{drive} & \text{if } s = \text{injured} \end{cases}$$

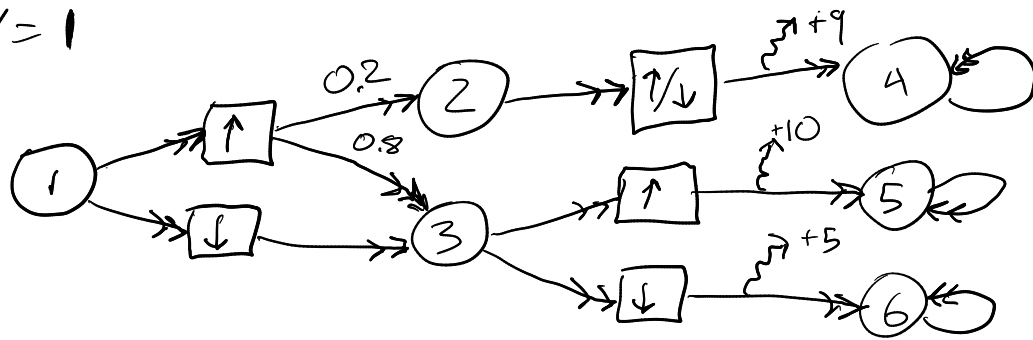
Guiding Questions

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

MDP Example: Up-Down Problem

$\gamma = 1$



Bellman Backup Algorithm

$U^*(s) \leftarrow 0$ for all terminal states

Repeat until all $U^*(s)$ are calculated:

find $U^*(s)$ for all states where
 $U^*(s)$ is known for all children

Extract $\pi^* = \operatorname{argmax}_a Q^*(s, a)$
 $= \operatorname{argmax}_a (R(s, a) + \gamma E[U^*(s')])$

$$U^*(s) = U^{\pi^*}(s)$$

expected sum of future rewards given
 that we follow the optimal policy

$$U^*(s) = \max_a (R(s, a) + \gamma E[U^*(s')])$$

$$\max_a (R(s, a) + \gamma \sum_{s'} T(s'|s, a) U^*(s'))$$

$$U^*(s) = \max_a Q^*(s, a) \quad Q^*(s, a)$$

s	a	$Q^*(s, a)$	$U^*(s)$
4			0
5			0
6			0
2	↑/↓	$R(2, \cdot) + (1 \cdot U^*(4))$ $9 + 0 = 9$	9
3	↑	$R(3, \uparrow) + (1 \cdot U^*(5)) = 10$	10
	↓	$R(3, \downarrow) + (1 \cdot U^*(6)) = 5$	
1	↑	$R(1, \uparrow) + (0.2 \cdot U^*(2) + 0.8 \cdot U^*(3))$ $0 + (0.2 \cdot 9 + 0.8 \cdot 10) = 9.8$	10
	↓	$R(1, \downarrow) + (1 \cdot U^*(3)) = 10$	

π^*

Dynamic Programming and Value Backup

Dynamic Programming and Value Backup

Bellman's Principle of Optimality: Every sub-policy in an optimal policy is locally optimal

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$

Policy Iteration

Algorithm: Policy Iteration



Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$
6. return π

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$  Policy Eval
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$  Policy Improvement Step
6. return π

(Policy iteration notebook)

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$
5. return U'

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$
5. return U'

- Returned U' will be close to U^* !

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_\infty > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U(s')) \quad \forall s \in S$
5. return U'

- Returned U' will be close to U^* !
- π^* is easy to extract: $\pi^*(s) = \arg \max (R(s, a) + \gamma E[U^*(s)])$

Bellman's Equations

Policy
Evaluation

$$U^\pi(s) = R(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim T(s'|s, \pi(s))} [U^\pi(s')]$$

Bellman's
Expectation Equation

Bellman Backup
Certificate of Optimality

$$U^*(s) = \max_a \left(R(s, a) + \gamma \mathbb{E}_{s' \sim T(s'|s, a)} [U^*(s')] \right)$$

Bellman's Optimality
Equation

Value Iteration

$$\begin{cases} U'(s) = \max_a \left(R(s, a) + \gamma \mathbb{E}_{s' \sim T(s'|s, a)} [U(s')] \right) \\ U'(s) = B[U](s) \end{cases}$$

Bellman Operator

VI:
initialize U, U'
while $\|U - U'\|_\infty > \epsilon$
 $U \leftarrow U'$
 $U' \leftarrow B[U]$

Guiding Questions

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

"In any small change he will have to consider only these quantitative indices (or "values") in which all the relevant information is concentrated; and by adjusting the quantities one by one, he can appropriately rearrange his dispositions without having to solve the whole puzzle ab initio, or without needing at any stage to survey it at once in all its ramifications."

-- F. A. Hayek, "The use of knowledge in society", 1945