

m

Markov Decision Processes

Last Time

- What does "Markov" mean in "Markov Process"?

$$s_{t+1} \perp s_{t-1}, \dots, s_0 \mid s_t \quad (\text{also } s_i \perp s_{t-1}, \dots, s_0 \mid s_t \quad \forall i > t)$$

$$p(s_{t+1} \mid s_t, \dots, s_0) = p(s_{t+1} \mid s_t)$$

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?

Decision Networks and MDPs




Decision Networks and MDPs

Decision Network

-  Chance node
-  Decision node
-  Utility node

Decision Networks and MDPs




Decision Network

-  Chance node
-  Decision node
-  Utility node

MDP Dynamic Decision Network

Decision Networks and MDPs

Decision Network




-  Chance node
-  Decision node
-  Utility node

MDP Dynamic Decision Network

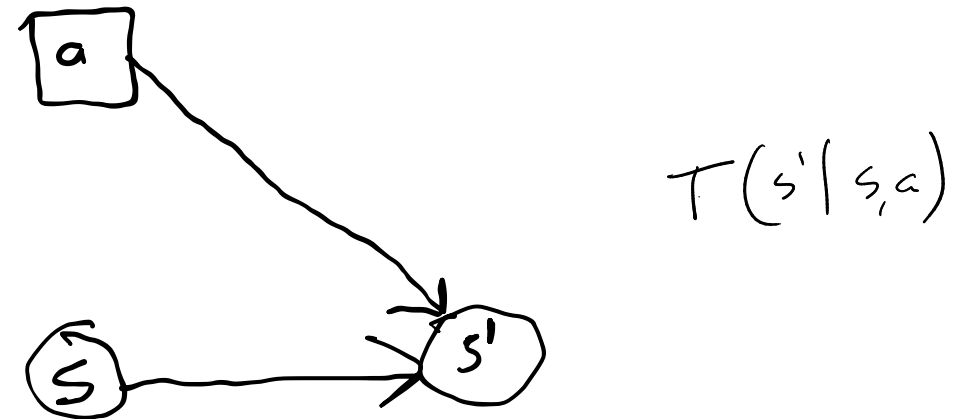


Decision Networks and MDPs

Decision Network

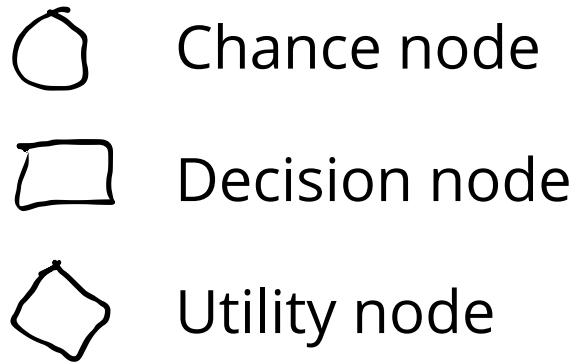
-  Chance node
-  Decision node
-  Utility node

MDP Dynamic Decision Network

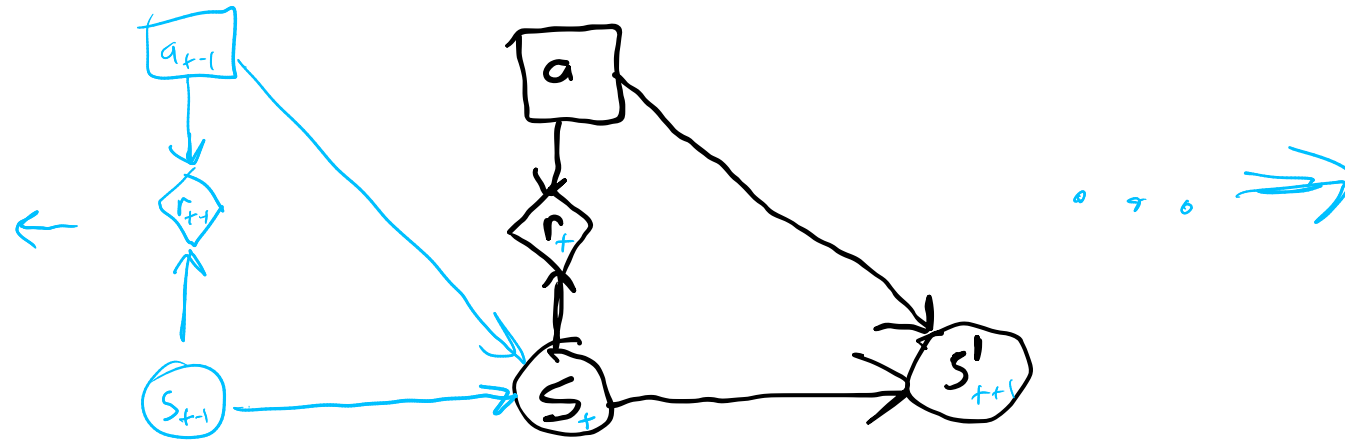


Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network



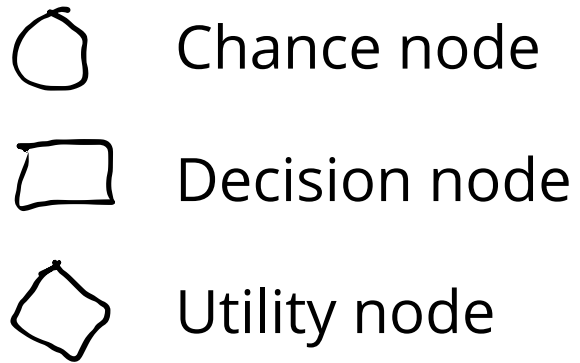
$$a = (a^1, a^2, a)$$

diagnostic

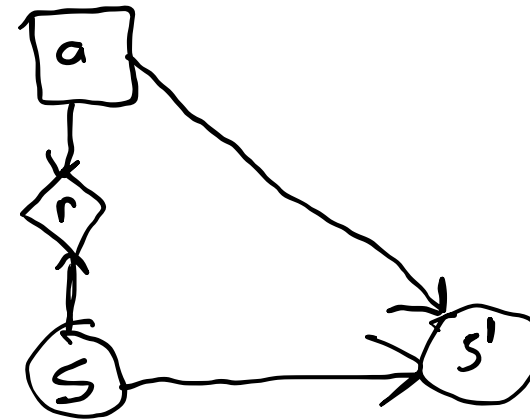
treatment

Decision Networks and MDPs

Decision Network



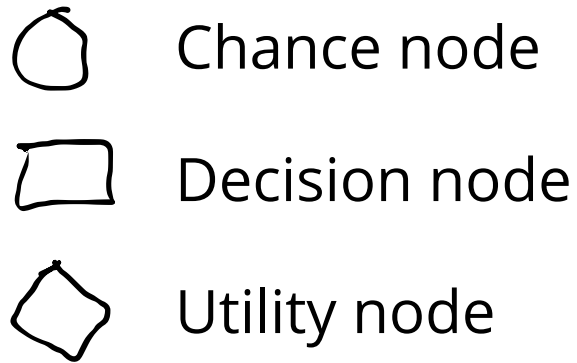
MDP Dynamic Decision Network



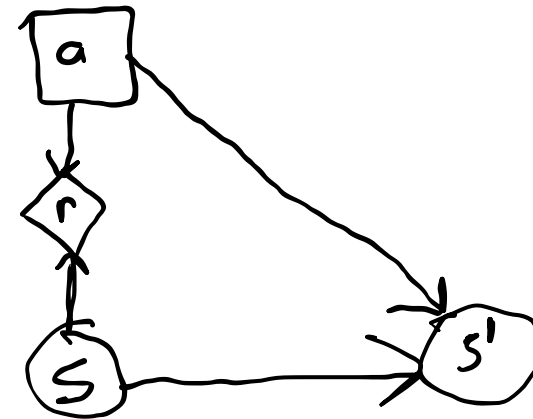
MDP Optimization problem

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network

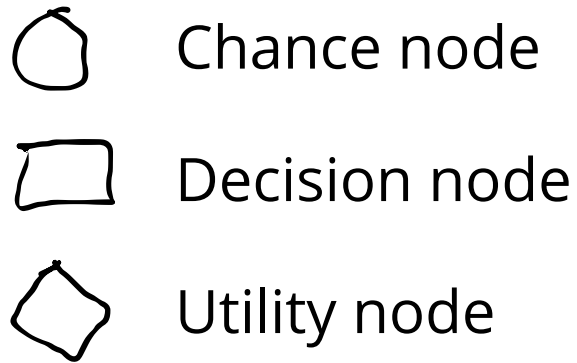


MDP Optimization problem

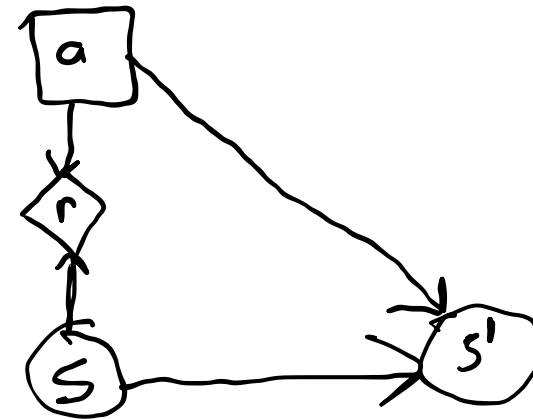
$$\text{maximize } \mathbb{E} \left[\sum_{t=0}^{\infty} r_t \right]$$

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network

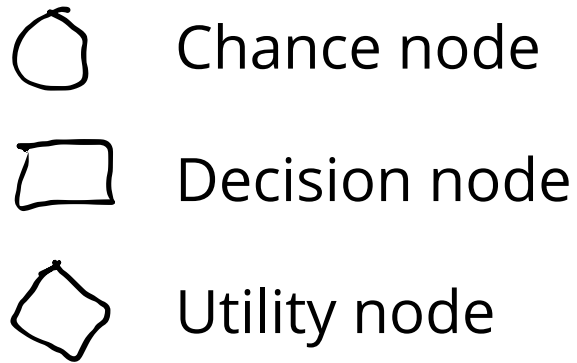


MDP Optimization problem

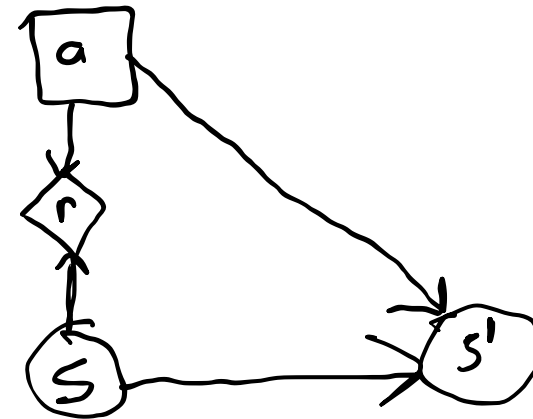
maximize $E \left[\sum_{t=1}^{\infty} r_t \right]$ Not well formulated!

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network



MDP Optimization problem

$$\text{maximize } \mathbb{E} \left[\sum_{t=1}^{\infty} r_t \right]$$

Not well formulated!
Infinite

Finite MDP Objectives

Finite MDP Objectives

1. Finite time

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if $\underline{r} \leq r_t \leq \bar{r}$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if $\underline{r} \leq r_t \leq \bar{r}$

then

$$\frac{\underline{r}}{1 - \gamma} \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1 - \gamma}$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if $\underline{r} \leq r_t \leq \bar{r}$

4. Terminal States

then

$$\frac{\underline{r}}{1 - \gamma} \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1 - \gamma}$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if $\underline{r} \leq r_t \leq \bar{r}$

4. Terminal States

Infinite time, but a terminal state (no reward, no leaving) is always reached with probability 1.

then

$$\frac{\underline{r}}{1 - \gamma} \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1 - \gamma}$$

Finite MDP Objectives

1. Finite time

$$\mathbb{E} \left[\sum_{t=0}^T r_t \right]$$

2. Average reward

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \sum_{t=0}^n r_t \right]$$

3. Discounting

$$\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$$

discount $\gamma \in [0, 1)$

typically 0.9, 0.95, 0.99

if $\underline{r} \leq r_t \leq \bar{r}$

4. Terminal States

Infinite time, but a terminal state (no reward, no leaving) is always reached with probability 1.

then

$$\frac{\underline{r}}{1 - \gamma} \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq \frac{\bar{r}}{1 - \gamma}$$

MDP "Tuple Definition"

MDP "Tuple Definition"

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

MDP "Tuple Definition"

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

MDP "Tuple Definition"

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

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

MDP "Tuple Definition"

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

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

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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 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 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 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 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 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 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 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 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

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

$\{1, 2, 3\}$
 \mathbb{R}^2
- A (action space) - set of all possible actions

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

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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

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

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

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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

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

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

MDP "Tuple Definition"

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

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

$\{1, 2, 3\}$
 $\{\text{healthy, pre-cancer, cancer}\}$

$(x, y) \in \mathbb{R}^2$
 $(s, i, r) \in \mathbb{N}^3$

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

$\{1, 2, 3\}$
 $\{\text{test, wait, treat}\}$

\mathbb{R}^2

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

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

$\{1, 2, 3\}$
 \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

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

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

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

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

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

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

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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 $\{0, 1\} \times \mathbb{R}^2$
 $\{\text{test, wait, treat}\}$
- T (transition distribution) - explicit or implicit ("generative")
model of how the state changes

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

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

$$\underline{s', r = G(s, a)}$$

MDP "Tuple Definition"

(S, A, T, R, γ) (and b 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

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

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

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

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

$\{\text{healthy, pre-cancer, cancer}\}$
 $(s, i, r) \in \mathbb{N}^3$
- A (action space) - set of all possible actions

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

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

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

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

$s', r = G(s, a)$
- b : initial state distribution

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

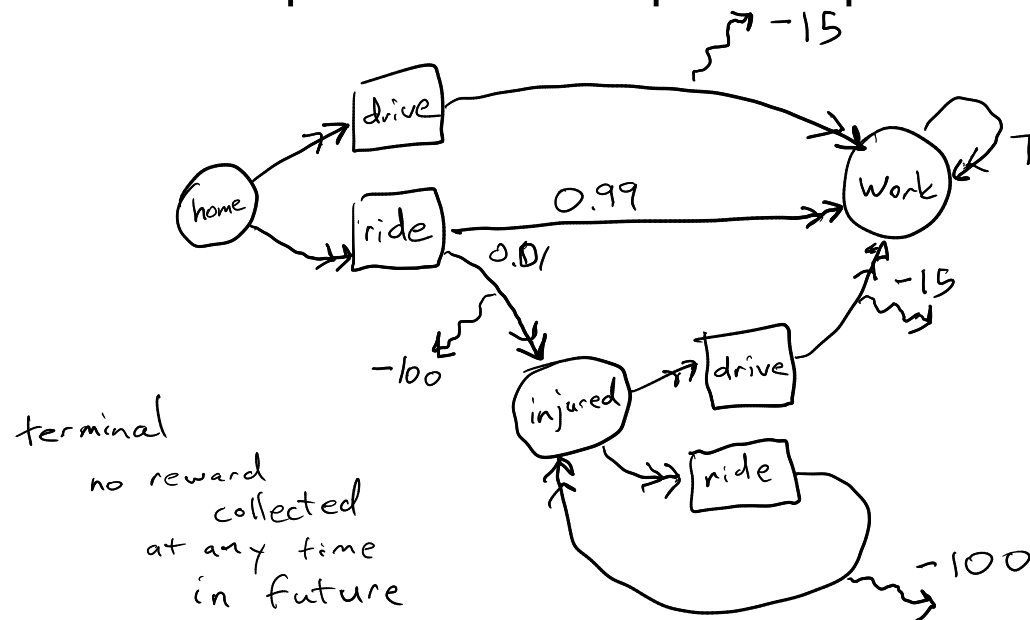
Imagine it's a cold day and you're ready to go to work. You have to decide whether to bike or drive.

- 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.



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

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

$$T^{\text{drive}} = \begin{matrix} & s' \\ s & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

$$T^{\text{ride}} = \begin{bmatrix} 0 & 0.99 & 0.01 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$R(s, a, s') = \begin{cases} -15 & \text{if } a = \text{drive} \\ -100 & \text{if } s' = \text{injured} \\ 0 & \text{o.w.} \end{cases} \quad \gamma = 0.99$$

Policies and Simulation

Policies and Simulation

- A *policy*, denoted with π , as in $a_t = \pi(s_t)$ is a function mapping every state to an action.
- When a policy is combined with a Markov decision process, it becomes a Markov stochastic process with

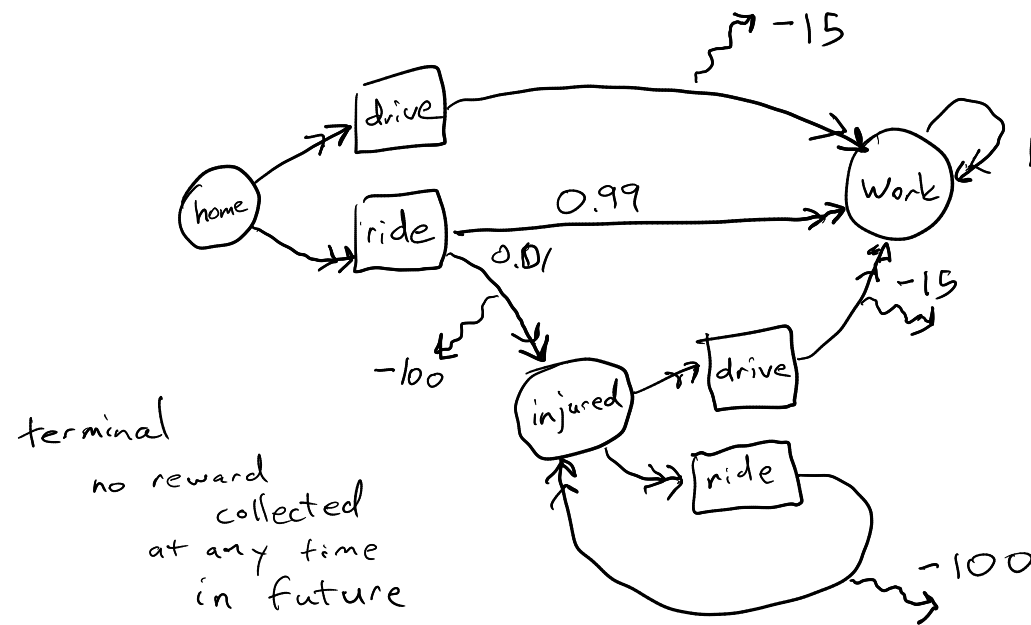
$$P(s' | s) = \underline{T(s' | s, \underline{\pi(s)})}$$

Simulate
Input: (S, A, T, R, γ, b) , π
Output: \hat{u} (accumulated reward)

$s \leftarrow \text{sample}(b)$
 $\hat{u} \leftarrow 0$
for t in $0 \dots T-1$ until $\gamma^t < \epsilon$
 $a \leftarrow \pi(s)$
 $s', r \leftarrow G(s, a)$
 $\hat{u} \leftarrow \hat{u} + \gamma^t r$
 $s \leftarrow s'$
return \hat{u}

Break

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



~~bike~~ drive -15

$$\text{bike: } 0.99 \cdot 0 + 0.01(-100 + -15) = -1.15$$

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

Policy Evaluation

$$U(\pi) = E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid a_t = \pi(s_t) \right] \quad r_t = R(s_t, a_t)$$

Naive:

$$U(\pi) = \sum_{t=0}^{\infty} \gamma^t \sum_{s_t \in S} P^{\pi}(s_t) R(s_t, \pi(s_t))$$

↙ marginal distribution of $s_t \mid \pi$

$$P^{\pi}(s_t) = \sum_{s_{t-1}} T(s_t \mid s_{t-1}, \pi(s_{t-1})) P^{\pi}(s_{t-1}) \quad \leftarrow$$

$$P^{\pi}(s_0) = b(s_0)$$

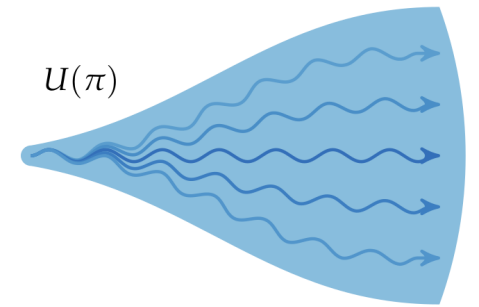
Monte Carlo Policy Evaluation

Monte Carlo Policy Evaluation

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

Monte Carlo Policy Evaluation

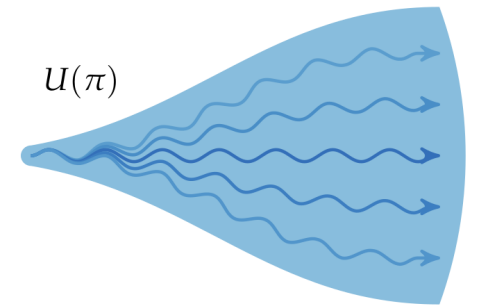
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



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

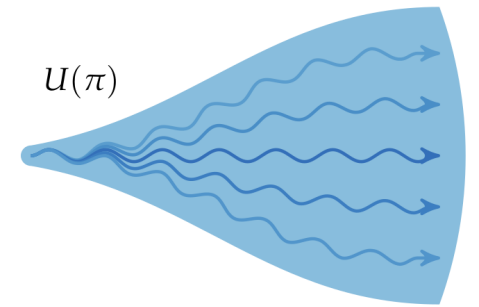


Monte Carlo Policy Evaluation

- 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

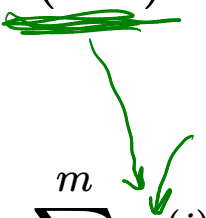
$$U(\pi) \approx \frac{1}{m} \sum_{i=1}^m R(\tau^{(i)})$$



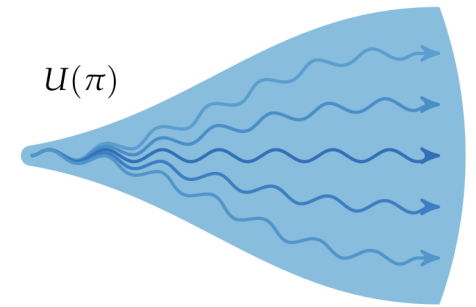
Monte Carlo Policy Evaluation

- 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

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


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



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

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

also a R.V.

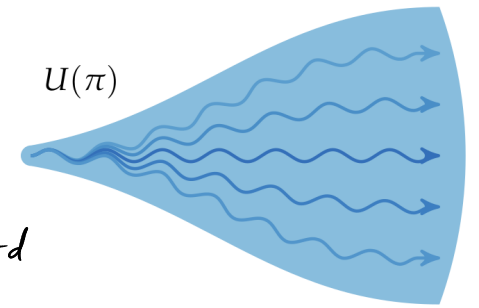
$$U(\pi) \approx \bar{u}_m = \frac{1}{m} \sum_{i=1}^m \hat{u}^{(i)}$$

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

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

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

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



Monte Carlo Policy Evaluation

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

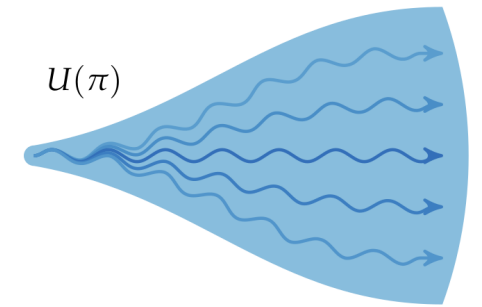
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

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

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



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

Monte Carlo Policy Evaluation

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

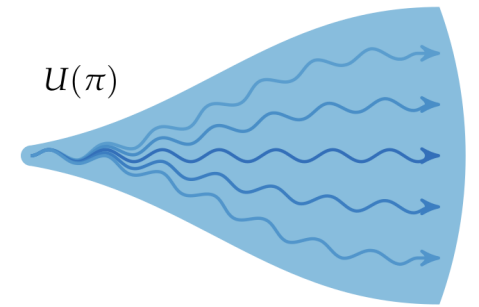
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

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

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



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

Value Function-Based Policy Evaluation

~~Discrete~~

Discrete MDPs only!

$$\begin{aligned}
 U^\pi(s) &= E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_t = \pi(s_t) \right] \\
 &= E[r_0 \mid s_0 = s, a_0 = \pi(s_0)] + E \left[\sum_{t=1}^{\infty} \gamma^t r_t \mid s_0 = s, a_t = \pi(s_t) \right] \\
 &= R(s, \pi(s)) + \sum_{s' \in S} T(s' | s, a) E \left[\sum_{t=1}^{\infty} \gamma^t r_t \mid s_0 = s, s_1 = s', a_t = \pi(s_t) \right]
 \end{aligned}$$

← Markov

$$\begin{aligned}
 &= R(s, \pi(s)) + \sum_{s' \in S} T(s' | s, a) \underbrace{E \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s', a_t = \pi(s_t) \right]}_{U^\pi(s')} \\
 &\quad \tau = t-1 \quad s_t \perp s_0 \mid s_1 \quad P(s_t | s_0, s_1) = P(s_t | s_1)
 \end{aligned}$$

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

$U^\pi(s')$

Bellman Expectation Eq.

$$\begin{aligned}
 \vec{U}^\pi &= \vec{U}^\pi \quad \vec{U}^\pi = \vec{R}^\pi + \gamma T^\pi \vec{U}^\pi \\
 \vec{R}^\pi[i] &= R(i, \pi(i)) \\
 T_{ij}^\pi &= T(j | i, \pi(i)) \\
 (I - \gamma T^\pi) \vec{U}^\pi &= \vec{R}^\pi
 \end{aligned}$$

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

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?