

HW 1 due tonight
Part way through!

Last Time

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

Stochastic Process $\{x_t\}$ is Markov if

$$P(x_{t+1} | x_t, x_{t-1}, \dots, x_0) = P(x_{t+1} | x_t)$$

W

Guiding Questions

Guiding Questions

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

Guiding Questions

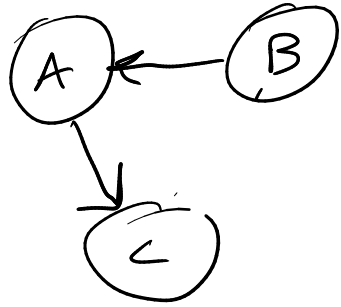
- What is a **Markov decision process**?
- What is a **policy**?

Guiding Questions

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

Decision Networks and MDPs

Bayes Net



$$P(C|A,B) = P(C|A)$$

Decision Networks and MDPs

Decision Network

Decision Networks and MDPs

Decision Network



Decision Networks and MDPs

Decision Network

 Chance node

Decision Networks and MDPs

Decision Network

 Chance node



Decision Networks and MDPs

Decision Network

 Chance node

 Decision node

Decision Networks and MDPs

Decision Network



Chance node



Decision node



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



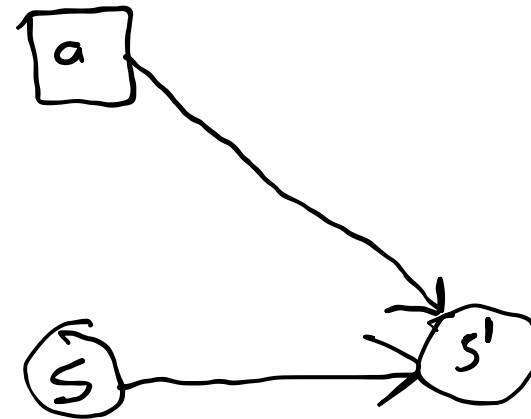
DBN for
a Markov
Stoch. Process

Decision Networks and MDPs

Decision Network

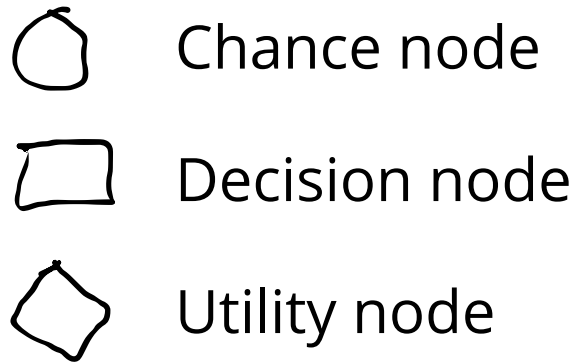
-  Chance node
-  Decision node
-  Utility node

MDP Dynamic Decision Network

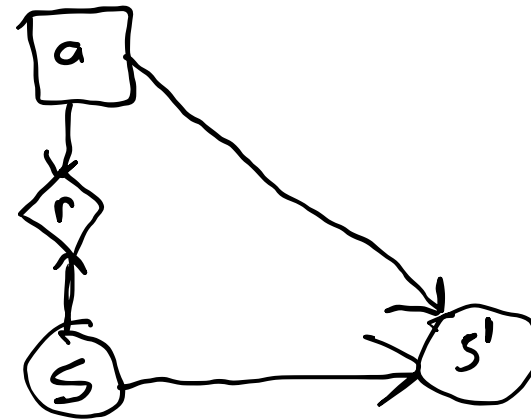


Decision Networks and MDPs

Decision Network

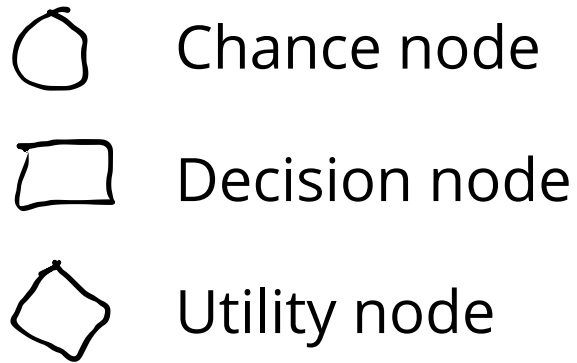


MDP Dynamic Decision Network

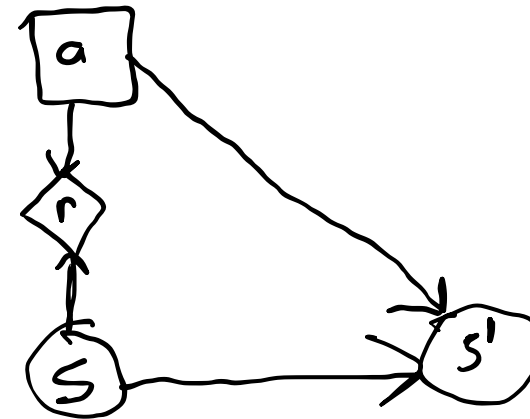


Decision Networks and MDPs

Decision Network



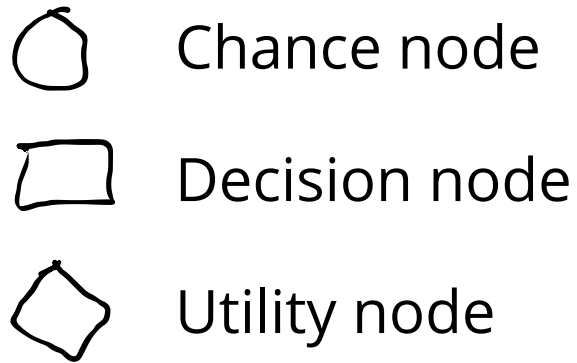
MDP Dynamic Decision Network



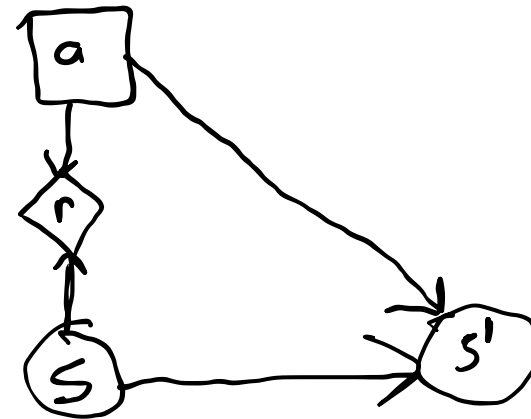
MDP Optimization problem

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network

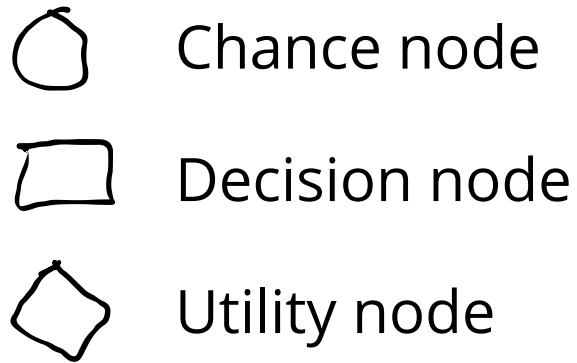


MDP Optimization problem

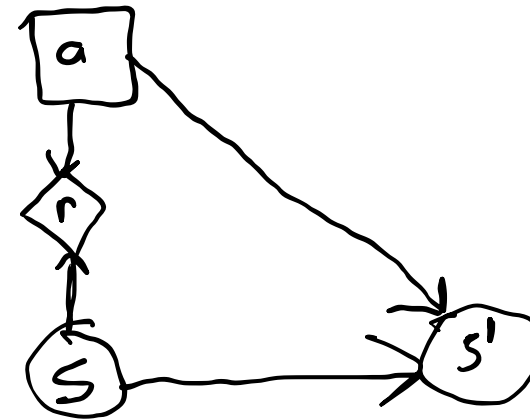
$$\text{maximize } \mathbb{E} \left[\sum_{t=1}^{\infty} r_t \right]$$

Decision Networks and MDPs

Decision Network



MDP Dynamic Decision Network



MDP Optimization problem

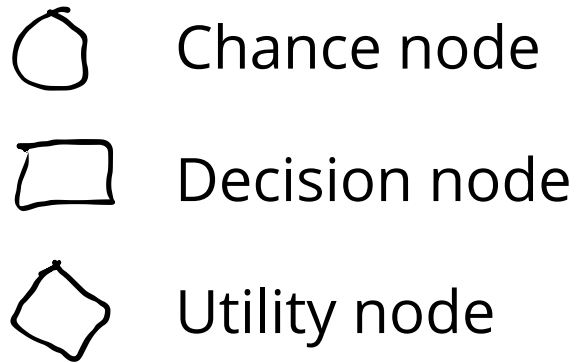
$$\text{maximize } \mathbb{E} \left[\sum_{t=1}^{\infty} r_t \right]$$

Not well formulated!

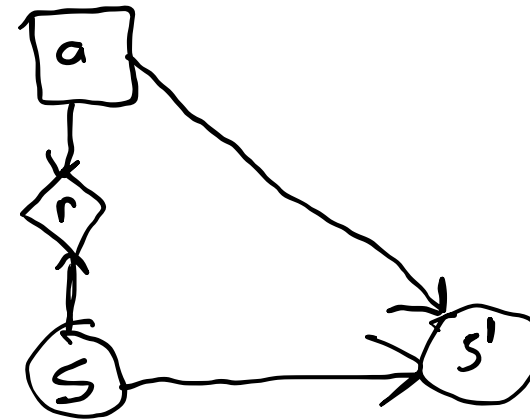
$$r_t = 1$$

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

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

Finite MDP Objectives

1. Finite time

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

2. Average reward

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

3. Discounting

Finite MDP Objectives

1. Finite time

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

2. Average reward

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

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

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

3. Discounting

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

discount $\gamma \in [0, 1)$
 $\gamma = 0.9$
 $\gamma^0 = 1$

Finite MDP Objectives

1. Finite time

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

2. Average reward

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

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

$$\gamma^0 = 1 \quad \gamma^1 = 0.99 \\ \gamma^2 = 0.99^2$$

Finite MDP Objectives

1. Finite time

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

2. Average reward

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

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

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

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

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

2. Average reward

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

3. Discounting

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

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

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

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

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

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

2. Average reward

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

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"

$$(\underset{\mathcal{S}}{S}, \underset{\mathcal{A}}{A}, T, R, \gamma)$$

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

$\{1, 2, 3\}$

$\{\text{healthy, pre-cancer, cancer}\}$

\mathbb{R}^2 $\{0, \overset{\downarrow}{1}\} \times \overset{\downarrow}{\mathbb{R}^4}$

MDP "Tuple Definition"

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

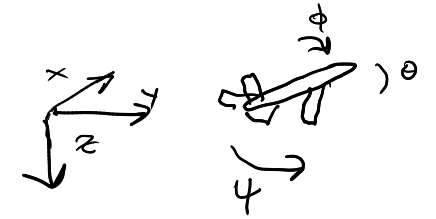
- S (state space) - set of all possible states

$\{1, 2, 3\}$ $(x, y) \in \mathbb{R}^2$ $\{0, 1\} \times \mathbb{R}^4$
 $\{\text{healthy, pre-cancer, cancer}\}$

state space

MDP "Tuple Definition"

$$(x, y, z, u, v, w, \phi, \theta, \psi, \rho, \alpha, r) \in \mathbb{R}^2$$



(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

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- 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
 $\{0, 1\} \times \mathbb{R}^2$

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

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

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

$(s, i, r) \in \mathbb{N}^3$

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

$\{1, 2, 3\}$
 $(x, y) \in \mathbb{R}^2$
 $\{0, 1\} \times \mathbb{R}^4$
- 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
 $\{0, 1\} \times \mathbb{R}^2$
- T (transition distribution) - explicit or implicit ("generative")
model of how the state changes

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

MDP "Tuple Definition"

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

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

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

$\{1, 2, 3\} \quad \mathbb{R}^2 \quad \{0, 1\} \times \mathbb{R}^2$
 $\{\text{test, wait, treat}\}$
- T (transition distribution) - explicit or implicit ("generative") model of how the state changes

$T(s' \mid s, a)$
 $\underbrace{s'}_{\mathbb{S}}, \underbrace{r}_{\mathbb{R}} = G(\underbrace{s}_{\mathbb{S}}, \underbrace{a}_{\mathbb{A}})$

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

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

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

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

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

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

$\{\text{test, wait, treat}\}$
 $T(s' \mid s, a)$
 $s', r = G(s, a)$
- R (reward function) - maps each state and action to a reward

$R(s, a)$

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

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

$s', r = G(s, a)$
- R (reward function) - maps each state and action to a reward

$R(s, a)$
- γ : discount factor

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

$T(s' | s, a) \leftarrow$
 $s', r = G(s, a) \leftarrow$
 $R(s, a)$
- R (reward function) - maps each state and action to a reward
- γ : discount factor
- b : initial state distribution

$$s_0 \sim b$$

$$S = \{1, 2, 3\} \quad A = \{1, 2\}$$

s'	$T(s' s=1, a=1)$
1	0.5
2	0.2
3	0.3

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, but you can't discover this until you start riding. After your crash, you limp home with pain equivalent to losing \$100.

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

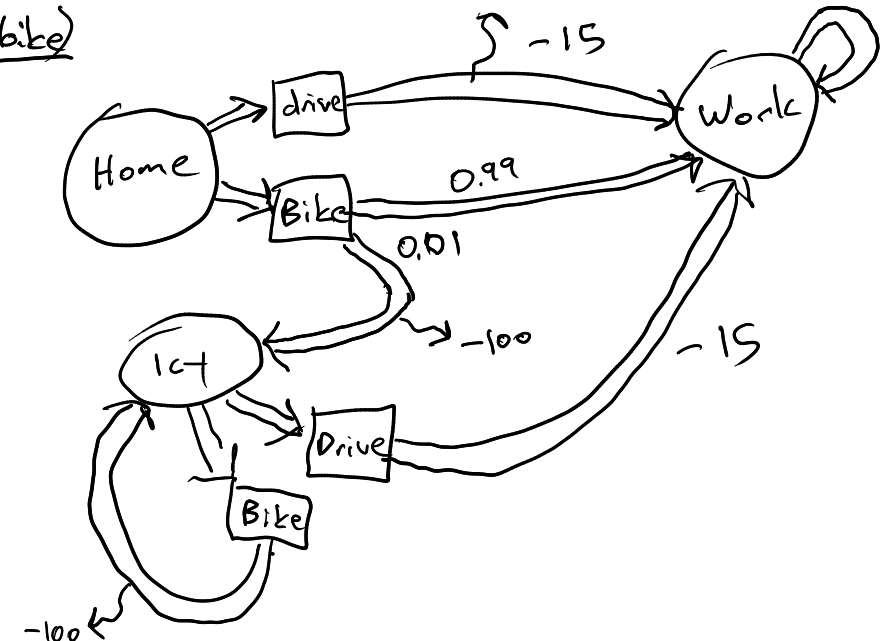
$$A = \{ \text{drive, bike} \}$$

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

$$\gamma = 0.99$$

$$T(s'|s,a) = \begin{cases} 1 & \text{if } s = \text{home, } a = \text{drive, } s' = \text{work} \\ 0.99 & \text{if } s = \text{home, } a = \text{bike, } s' = \text{work} \\ \vdots & \\ 0 & \text{if } a = \text{drive, } s' = \text{icy} \end{cases}$$

s'	$T(s' s=\text{home}, a=\text{bike})$
work	0.99
icy	0.01
home	0



Policy

Policy

- A *policy*, denoted with, π as in $a_t = \pi(s_t)$ is a function mapping every state to an action.

Policy

- A *policy*, denoted with, π as in $a_t = \pi(s_t)$ is a function mapping every state to an action.
- An optimal policy, π^* , maximizes the sum of expected rewards:

$$\pi^* = \underset{\pi}{\operatorname{argmax}} \operatorname{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \right]$$

$$a_t = \pi(s_t)$$

Breakout Rooms

- Name, "I think Colorado winters are _____"
 - Suggest a policy that you think is optimal for the icy day problem
-

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

MDP Simulation

MDP Simulation

Algorithm: Rollout Simulation

Given: MDP (S, A, R, T, γ, b)

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

$$\gamma^T < \epsilon$$

$\hat{u} \leftarrow 0$

for t in $0 \dots T - 1$

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

$$\begin{aligned} &T(s'|s, a) \\ &R(s, a) \end{aligned}$$

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

$s \leftarrow s'$

return \hat{u}

Policy Evaluation

Policy Evaluation

- *Evaluating* a policy means determining the expected reward.

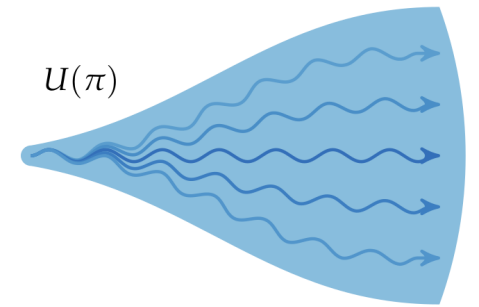
.

Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

Policy Evaluation

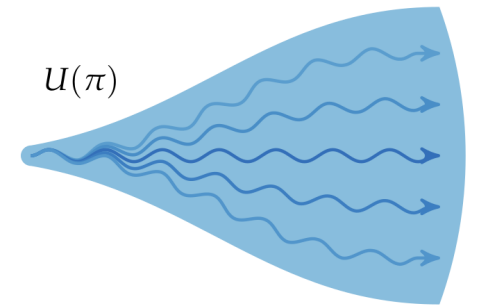
- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*



Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- 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

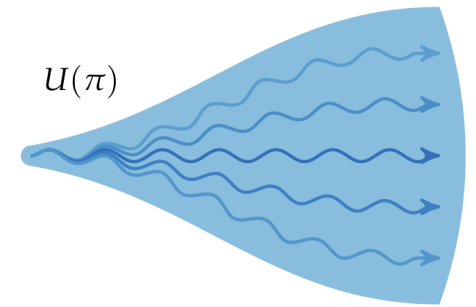


Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- 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

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



Policy Evaluation

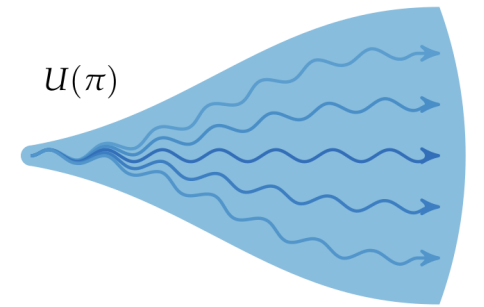
- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

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

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

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

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



Policy Evaluation

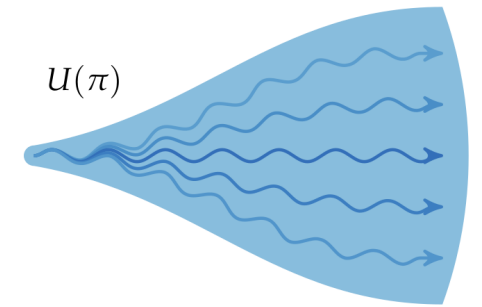
- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

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

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

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

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



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

Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

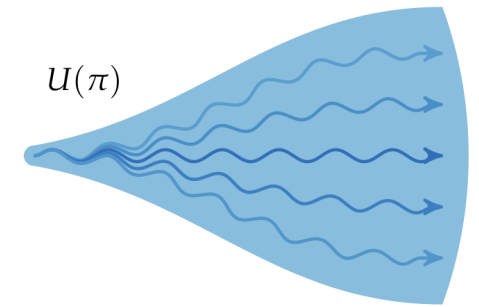
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

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

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 ?

Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- 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

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

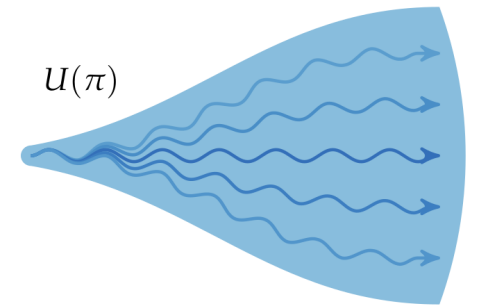
also an R.V.

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

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

C.L.T.

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



Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

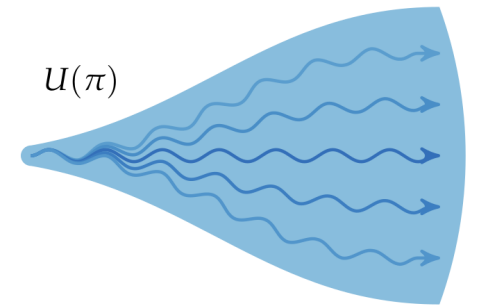
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

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

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 ?

C.L.T. $\frac{\bar{u}_m - U(\pi)}{\sigma_m / \sqrt{m}} \xrightarrow{d} \mathcal{N}(0, 1)$

Policy Evaluation

- *Evaluating* a policy means determining the expected reward.
- Running a large number of simulations and averaging the accumulated reward is called *Monte Carlo Evaluation*

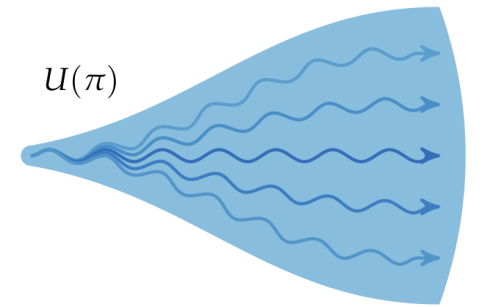
Let $\tau = (s_0, a_0, r_0, s_1, \dots, s_T)$ be a *trajectory* of the MDP

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

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 ?

C.L.T. $\frac{\bar{u}_m - U(\pi)}{\sigma_m / \sqrt{m}} \xrightarrow{d} \mathcal{N}(0, 1)$

$$\text{s.e.m.} = \frac{\text{std}(\hat{u})}{\sqrt{m}}$$

Guiding Questions

Guiding Questions

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

(S, A, T, R, γ)

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?