

Markov Decision Process

Announcements

TA office hours are posted

HW0 is due Wednesday

Programming assignment 1 will be out on Wednesday

Reading Materials: Reinforcement Learning: Theory & Algorithms

<https://rltheorybook.github.io/>

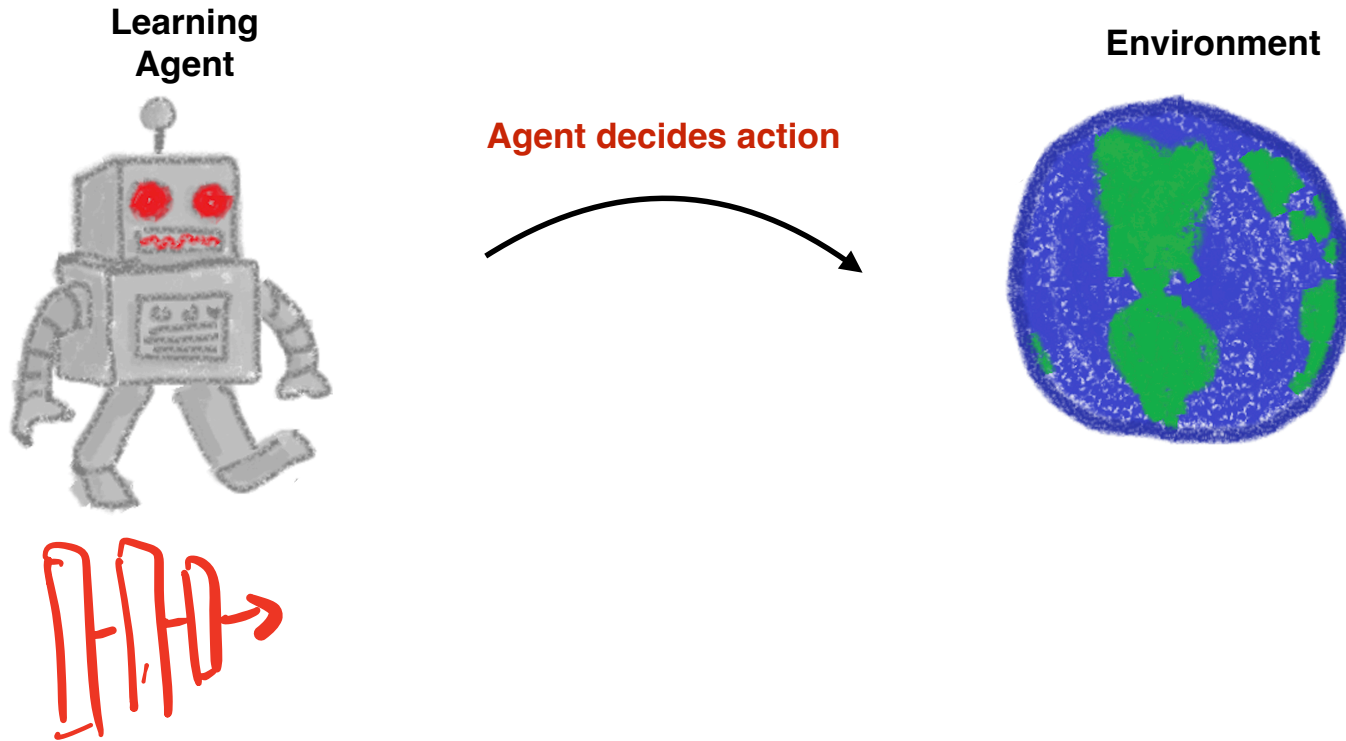
This is an extremely advanced RL book, so we will pick **specific subsections** for you to read

Please let us know if you find any typos or mistakes in the book

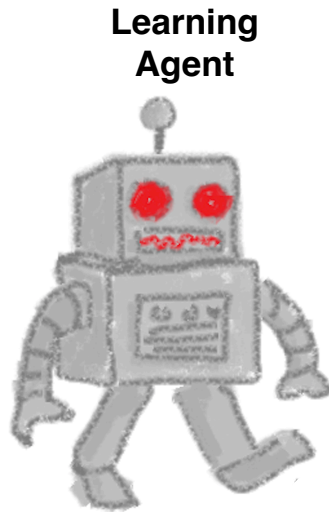
Outlines:

1. Definitions of Markov Decision Process
2. Value functions (V and Q functions)
3. Bellman equations

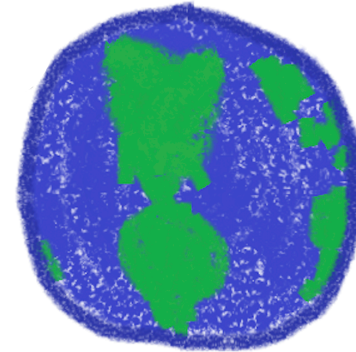
The Mathematical framework: Infinite horizon **Markov Decision Process**



The Mathematical framework: Infinite horizon **Markov Decision Process**



Environment



Agent decides action



Send **reward** and **next state** from a Markovian transition dynamics

$$r(s, a), s' \sim P(\cdot | s, a)$$

$P(\cdot | s, a)$: distribution over the next state

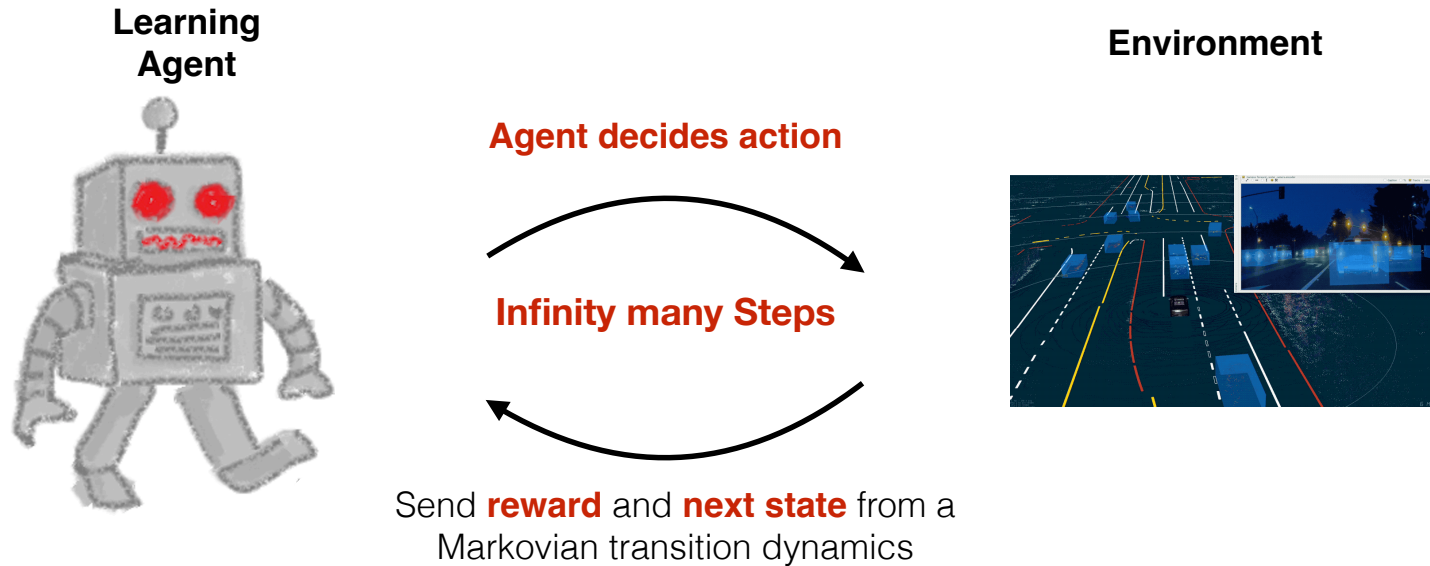
The Mathematical framework: Infinite horizon **Markov Decision Process**



$$r(s, a), s' \sim P(\cdot | s, a)$$

$P(\cdot | s, a)$: distribution over the next state

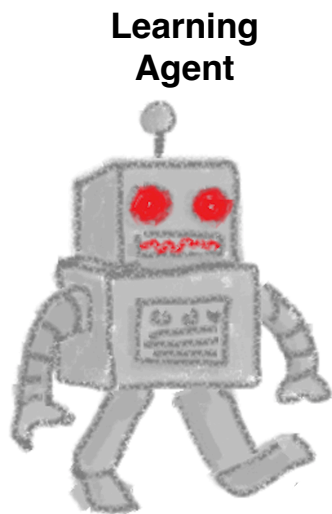
The Mathematical framework: Infinite horizon **Markov Decision Process**



$$r(s, a), s' \sim P(\cdot | s, a)$$

$P(\cdot | s, a)$: distribution over the next state

The Mathematical framework: Infinite horizon **Markov Decision Process**



Agent decides action



Infinity many Steps



Environment

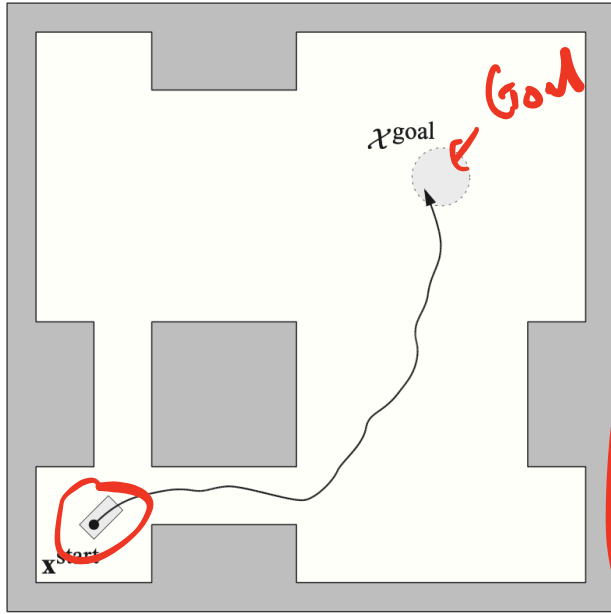


Send **reward** and **next state** from a Markovian transition dynamics

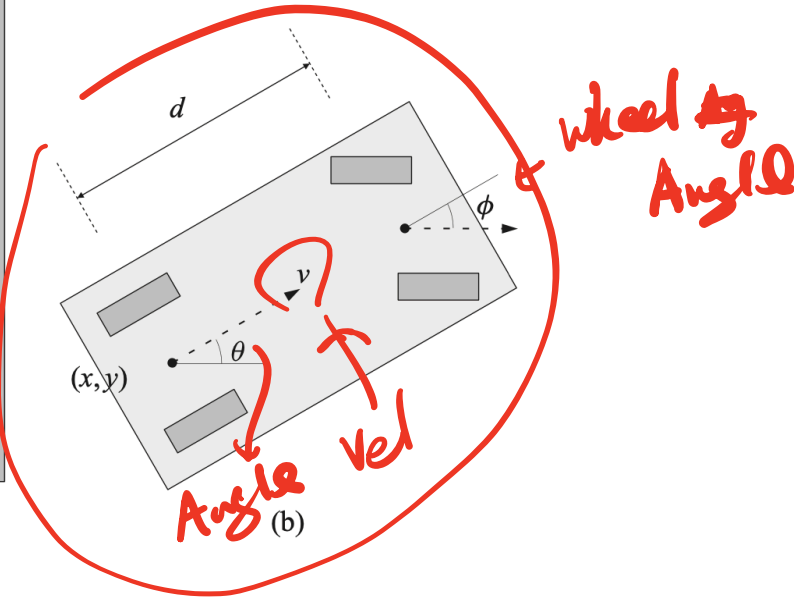
$$r(s, a), s' \sim P(\cdot | s, a)$$

$P(\cdot | s, a)$: distribution over the next state

Example: 2-D simple car navigation



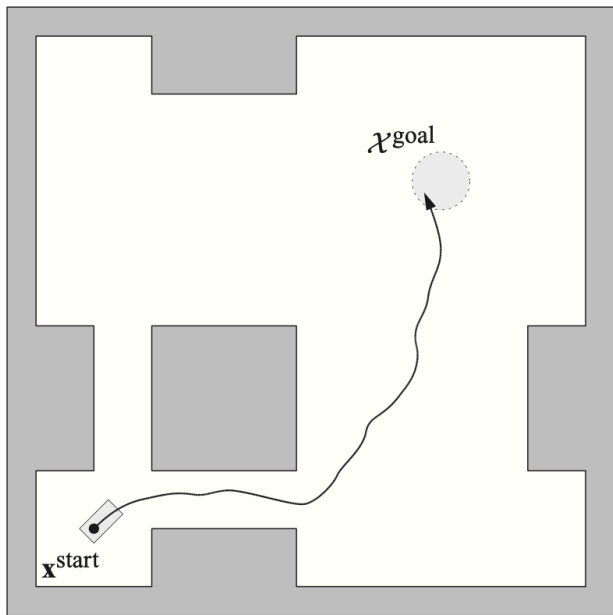
(a)



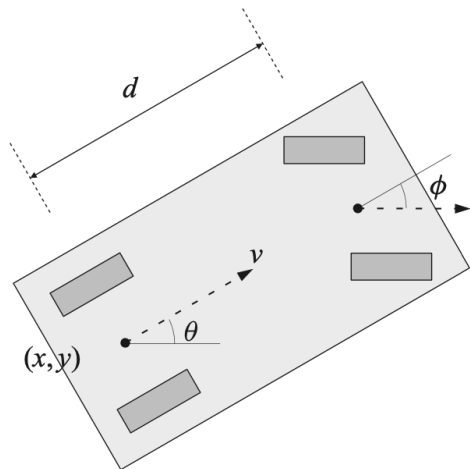
(b)

Example: 2-D simple car navigation

$$s = [x, y, \theta, v]^T \in \mathbb{R}^4$$

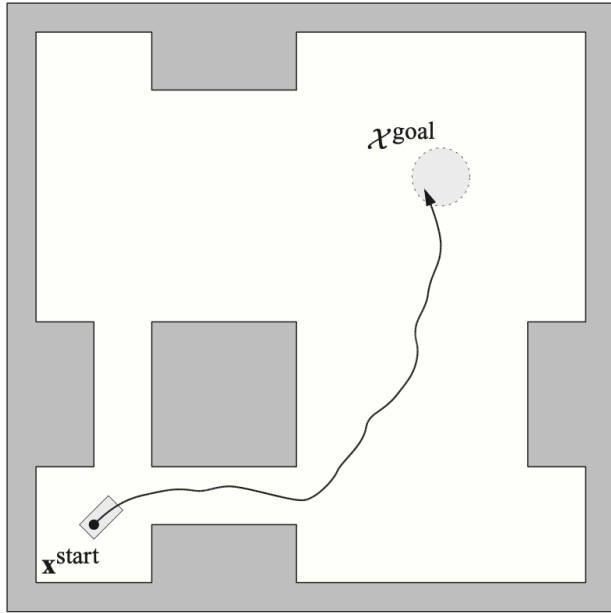


(a)

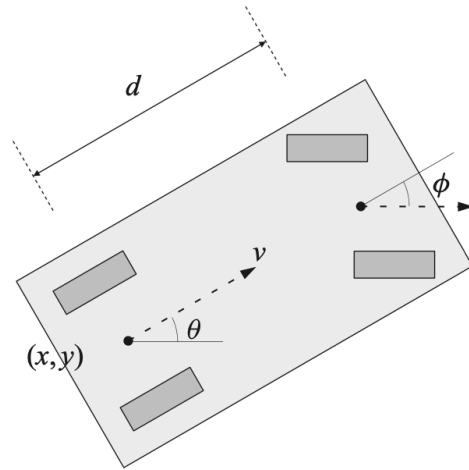


(b)

Example: 2-D simple car navigation



(a)



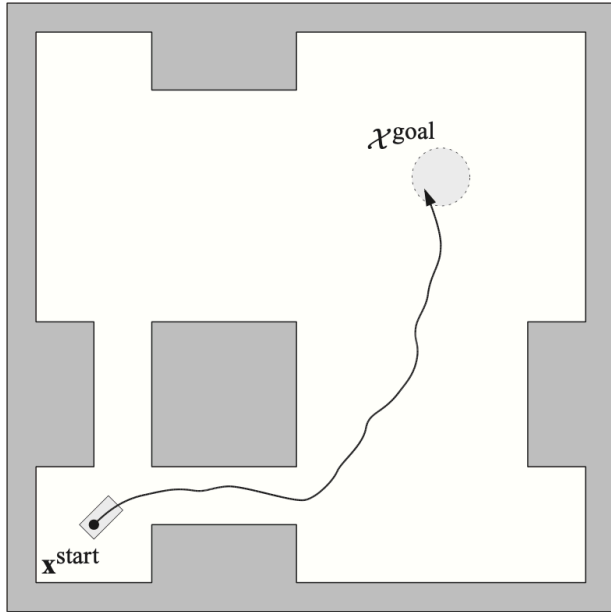
(b)

$$s = [x, y, \theta, v]^T \in \mathbb{R}^4$$

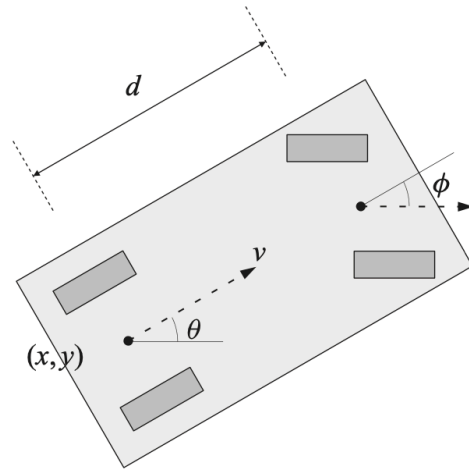
$$a = [\alpha, \phi]^T \in \mathbb{R}^2$$

\uparrow
Acc \uparrow wheel Angle

Example: 2-D simple car navigation



(a)



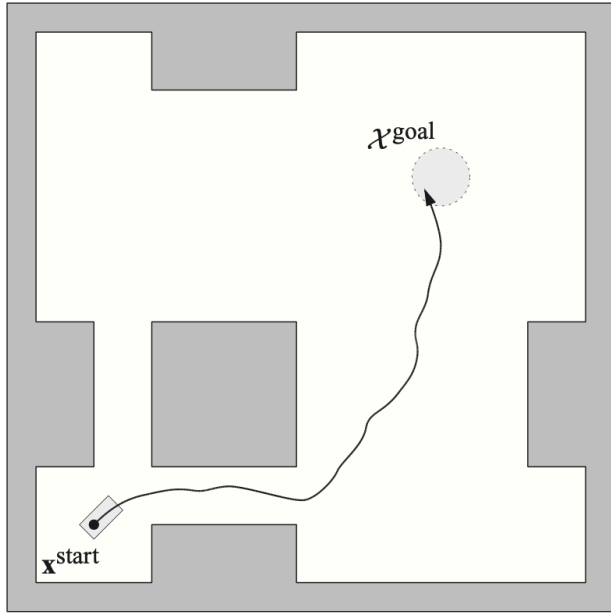
(b)

$$s = [x, y, \theta, v]^T \in \mathbb{R}^4$$

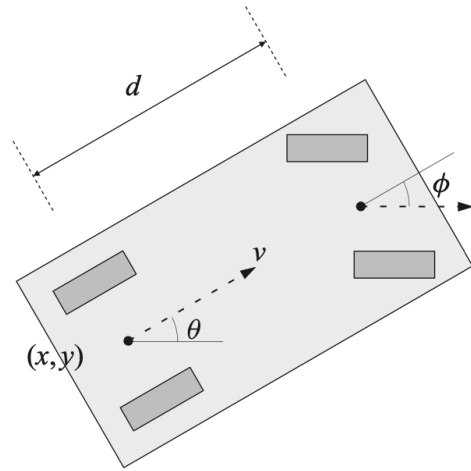
$$a = [\alpha, \phi]^T \in \mathbb{R}^2$$

$$r(s, a) = \begin{cases} 100 & (x, y) \in \mathcal{X}_{\text{goal}} \\ -1 & \text{hit obstacles} \\ 0 & \text{else} \end{cases}$$

Example: 2-D simple car navigation



(a)



(b)

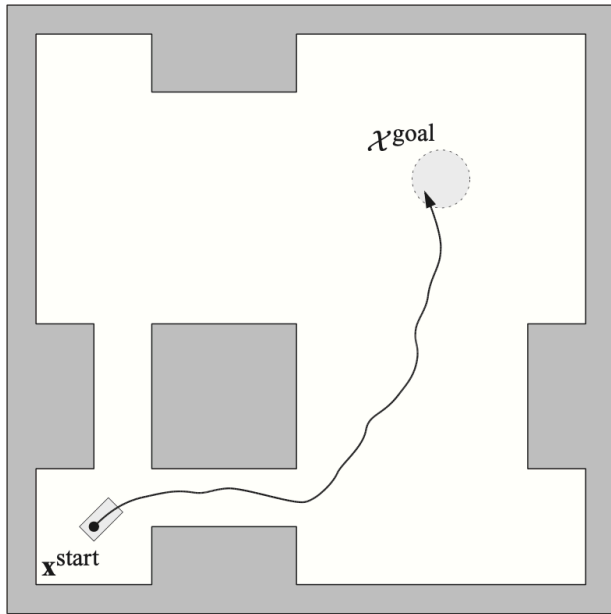
$$s = [x, y, \theta, v]^T \in \mathbb{R}^4$$

$$a = [\alpha, \phi]^T \in \mathbb{R}^2$$

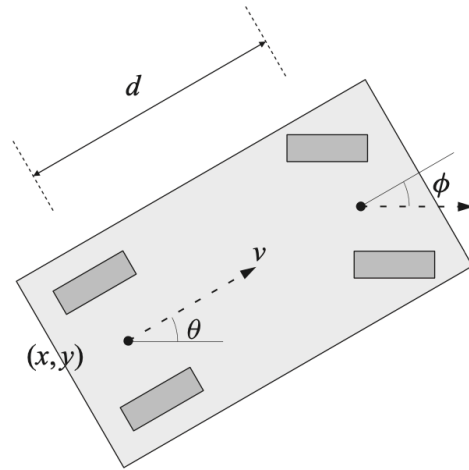
$$r(s, a) = \begin{cases} 100 & (x, y) \in \mathcal{X}_{\text{goal}} \\ -1 & \text{hit obstacles} \\ 0 & \text{else} \end{cases}$$

$$s' = f(s, a) + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, I)$$

Example: 2-D simple car navigation



(a)



(b)

$$s = [x, y, \theta, v]^T \in \mathbb{R}^4$$

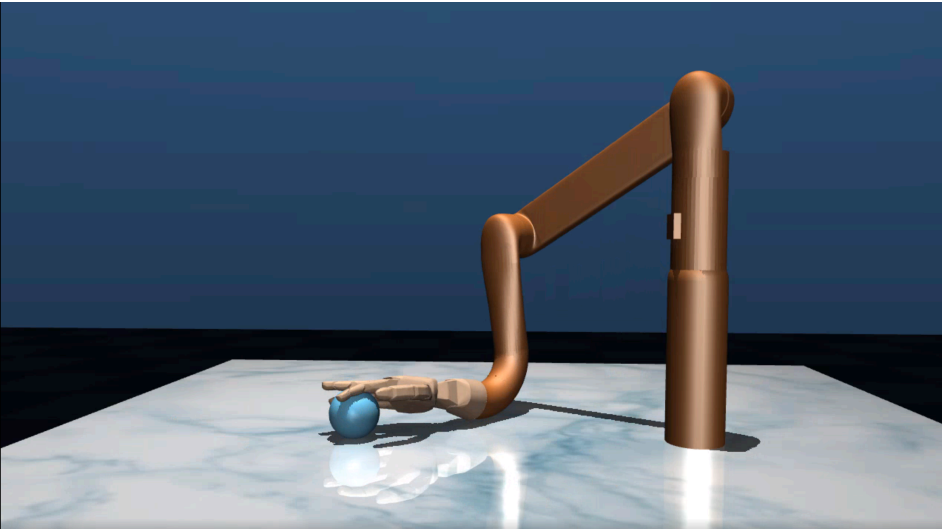
$$a = [\alpha, \phi]^T \in \mathbb{R}^2$$

$$r(s, a) = \begin{cases} 100 & (x, y) \in \mathcal{X}_{\text{goal}} \\ -1 & \text{hit obstacles} \\ 0 & \text{else} \end{cases}$$

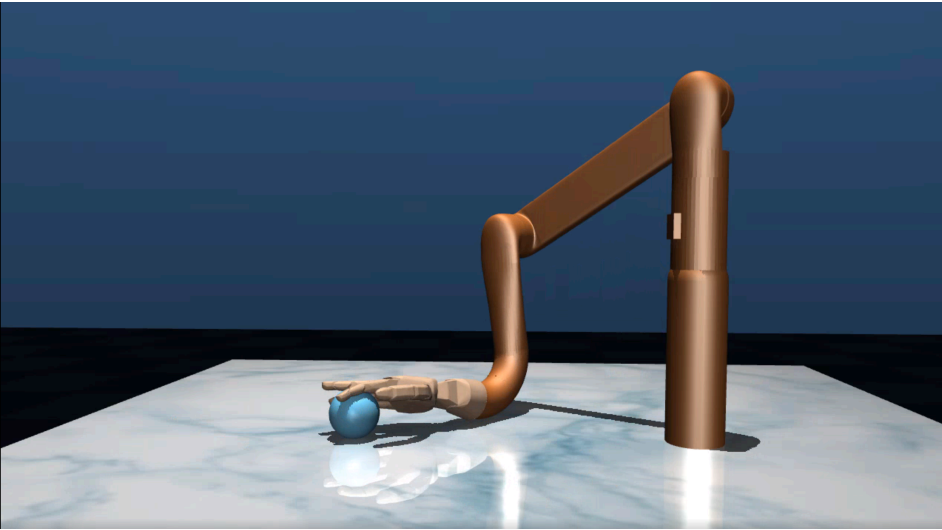
$$s' = f(s, a) + \epsilon, \text{ where } \epsilon \sim \mathcal{N}(0, I)$$

$$f(s, a) = \begin{bmatrix} x + \tau v \cos \theta \\ y + \tau v \sin \theta \\ \theta + \tau v \tan(\phi)/d \\ v + \tau \alpha \end{bmatrix}$$

Example:
robot hand needs to pick the ball and hold it in a goal (x,y,z) position



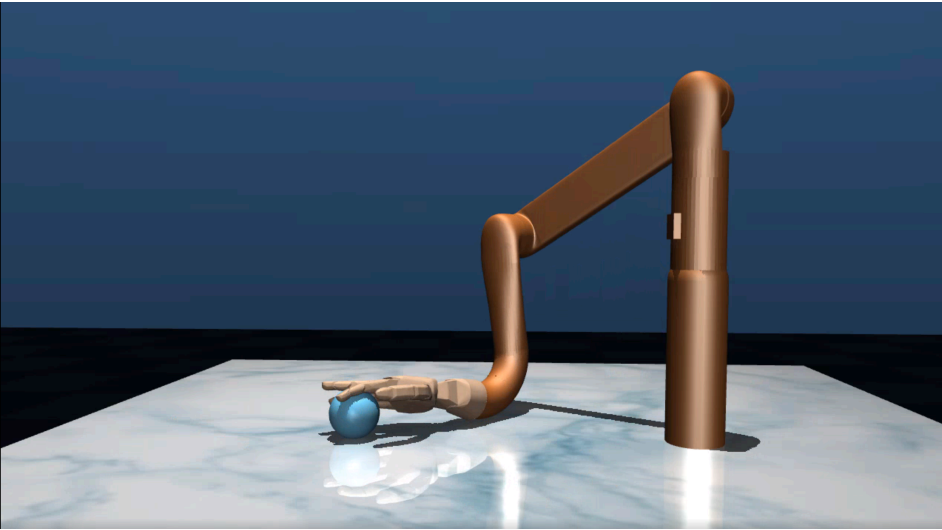
Example:
robot hand needs to pick the ball and hold it in a goal (x,y,z) position



State s : robot configuration (e.g., joint angles)
and the ball's position

Example:

robot hand needs to pick the ball and hold it in a goal (x,y,z) position

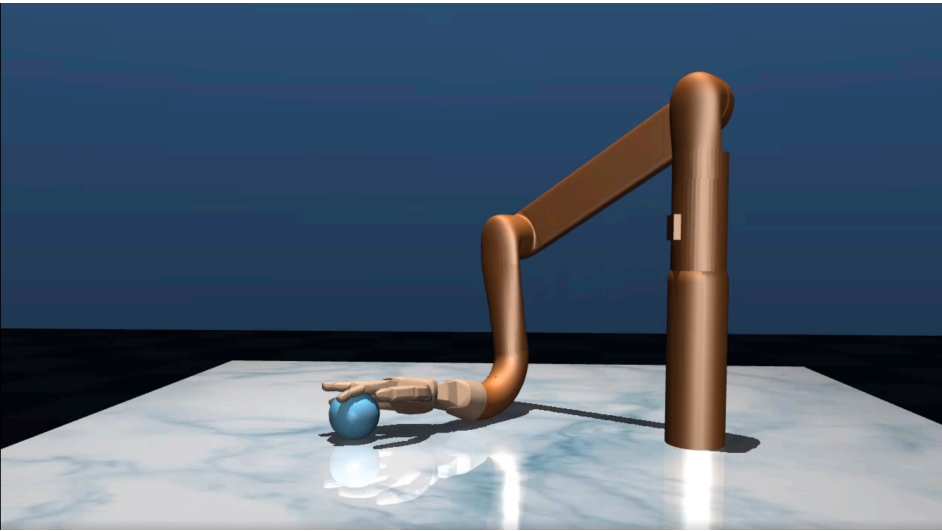


State s : robot configuration (e.g., joint angles) and the ball's position

Action a : Torque on joints in arm & fingers

Example:

robot hand needs to pick the ball and hold it in a goal (x,y,z) position



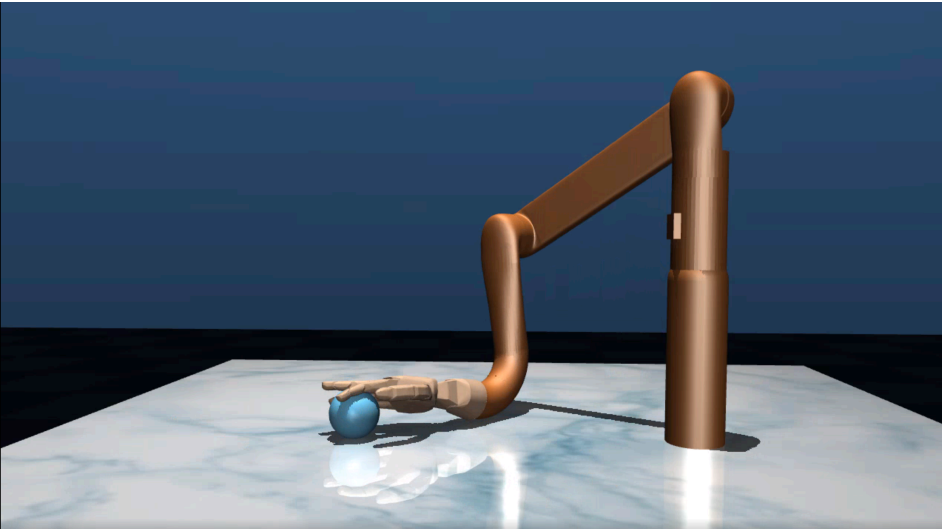
State s : robot configuration (e.g., joint angles) and the ball's position

Action a : Torque on joints in arm & fingers

Transition $s' \sim P(\cdot | s, a)$: physics + some noise

Example:

robot hand needs to pick the ball and hold it in a goal (x,y,z) position



State s : robot configuration (e.g., joint angles) and the ball's position

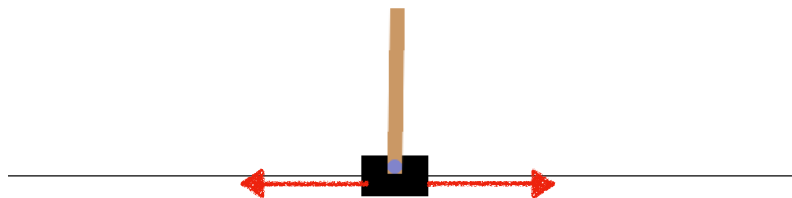
Action a : Torque on joints in arm & fingers

Transition $s' \sim P(\cdot | s, a)$: physics + some noise

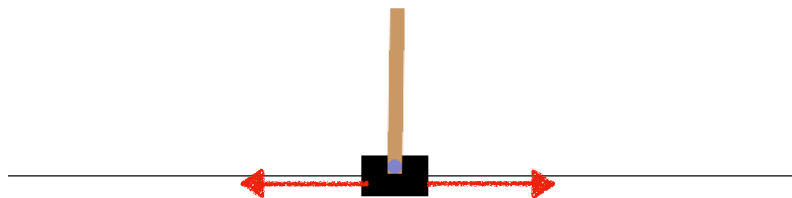
Cost $c(s, a)$: torque magnitude + dist to goal

Example: OpenAI Gym demonstrations

Cart pole



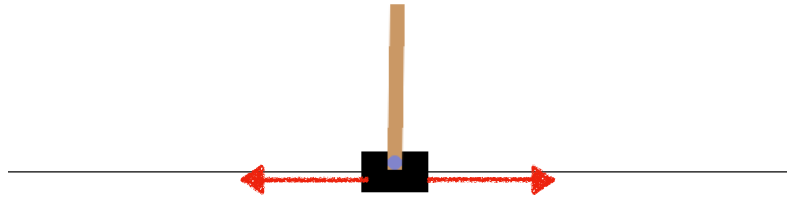
Example: OpenAI Gym demonstrations



Example: OpenAI Gym demonstrations

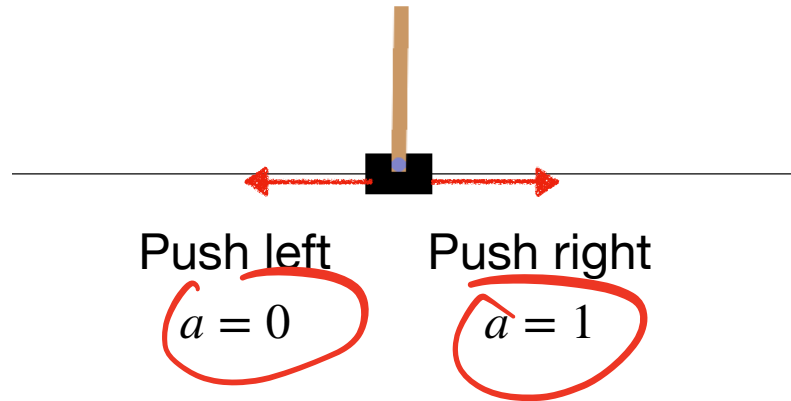
State = [cart pos, cart velocity,
pole angle, pole angular velocity]

$\in \mathbb{R}^4$
=



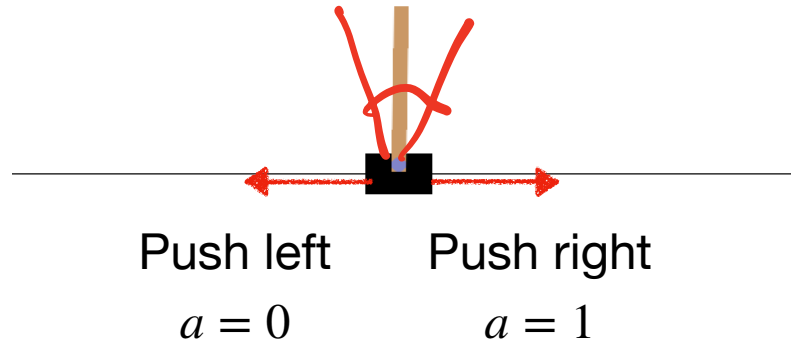
Example: OpenAI Gym demonstrations

State = [cart pos, cart velocity,
pole angle, pole angular velocity]



Example: OpenAI Gym demonstrations

State = [cart pos, cart velocity,
pole angle, pole angular velocity]



$$r(s, a) = \begin{cases} 1 & \text{pole angle} \in [-12^\circ, 12^\circ], \\ 0 & \text{else} \end{cases}$$

Policy

$$\mathcal{S} \rightarrow \mathcal{A}$$

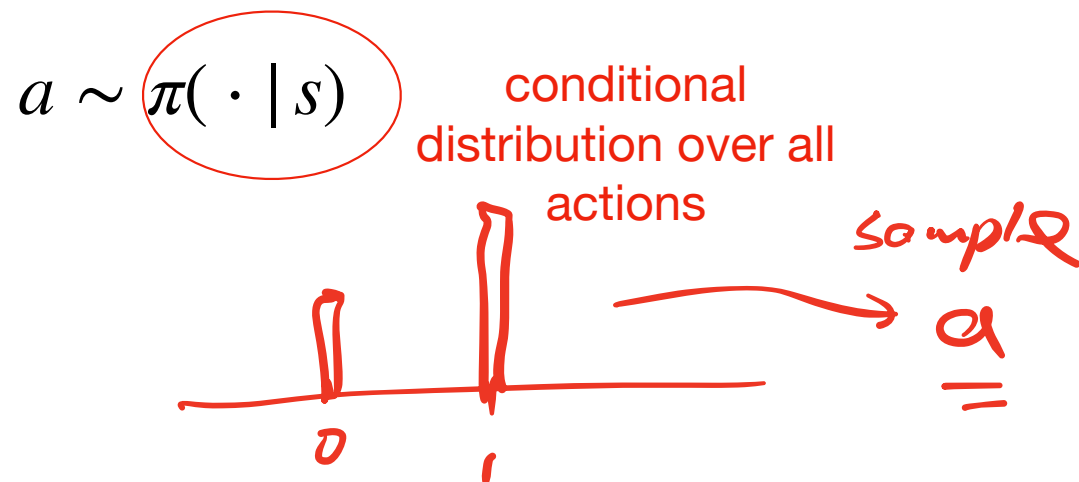
A mapping from state to action (what action should I take if I'm in this state...)

$$a \sim \pi(\cdot | s)$$

Policy

$$\mathcal{S} \rightarrow \mathcal{A}$$

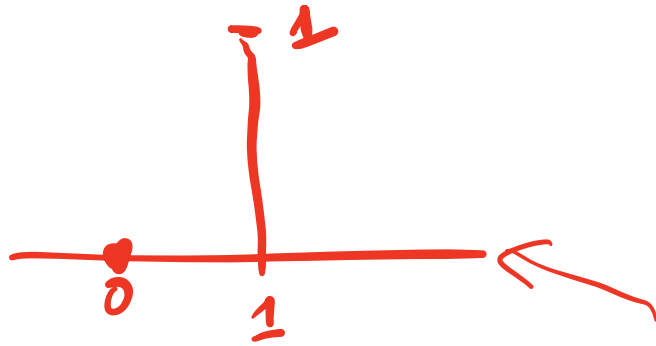
A mapping from state to action (what action should I take if I'm in this state...)



Policy

$$\mathcal{S} \rightarrow \mathcal{A}$$

A mapping from state to action (what action should I take if I'm in this state...)



$$a \sim \pi(\cdot | s)$$

conditional
distribution over all
actions

Deterministic vs stochastic?

Policy

$$\mathcal{S} \rightarrow \mathcal{A}$$

A mapping from state to action (what action should I take if I'm in this state...)

$$a \sim \pi(\cdot | s)$$

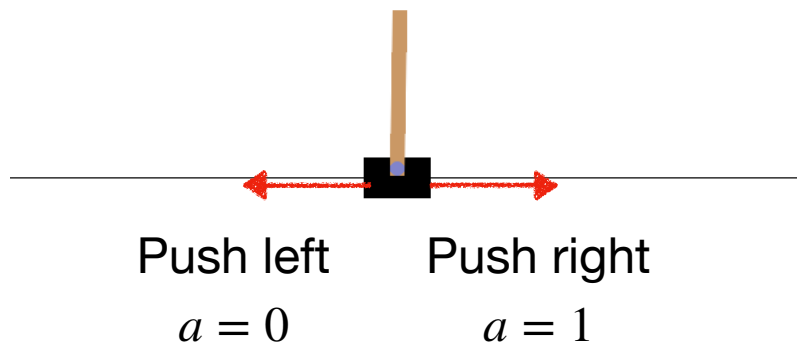
conditional
distribution over all
actions

Deterministic vs stochastic?

Q: Assume \mathcal{S} state and \mathcal{A} actions, how many different deterministic policies we can have?

Example: OpenAI Gym demonstrations

State = [cart pos, cart velocity,
pole angle, pole angular velocity]

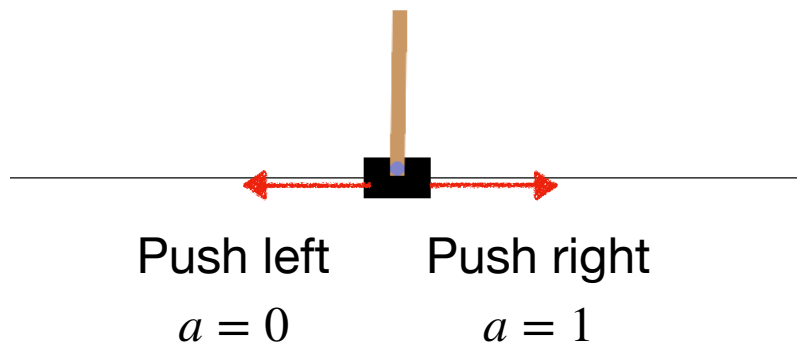


$$r(s, a) = \begin{cases} 1 & \text{pole angle} \in [-12^\circ, 12^\circ], \\ 0 & \text{else} \end{cases}$$

Policy 1: uniform random
 $\pi(0 | s) = \pi(1 | s) = 0.5, \forall s$

Example: OpenAI Gym demonstrations

State = [cart pos, cart velocity,
pole angle, pole angular velocity]

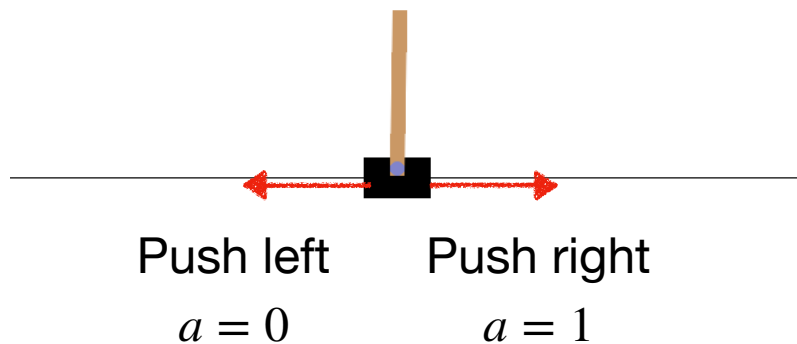


$$r(s, a) = \begin{cases} 1 & \text{pole angle} \in [-12^\circ, 12^\circ], \\ 0 & \text{else} \end{cases}$$

Policy 1: uniform random
 $\pi(0 | s) = \pi(1 | s) = 0.5, \forall s$

Example: OpenAI Gym demonstrations

State = [cart pos, cart velocity,
pole angle, pole angular velocity]




$$r(s, a) = \begin{cases} 1 & \text{pole angle} \in [-12^\circ, 12^\circ], \\ 0 & \text{else} \end{cases}$$

Policy 1: uniform random
 $\pi(0 | s) = \pi(1 | s) = 0.5, \forall s$

Policy 2: adaptive

$$\pi(s) = \begin{cases} 0(\text{left}) & \text{if pole angle} < 0 \\ 1(\text{right}) & \text{else} \end{cases}$$

Outlines:

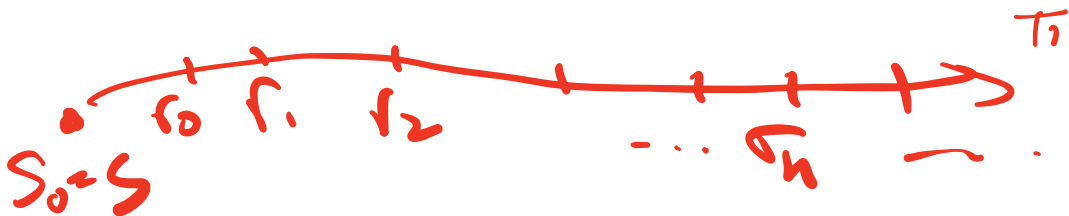
-  1. Definitions of Markov Decision Process
2. Value functions (V and Q functions)
3. Bellman equations

Time step Performance of a policy π

Expected total reward of a policy π :

$$V^\pi(s) = \mathbb{E} \left[\underbrace{r_0}_{\Delta} + \gamma \underbrace{r_1}_{\Delta} + \gamma^2 \underbrace{r_2}_{\Delta} + \dots + \gamma^h r_h + \dots \mid \underbrace{s_0 = s, \pi} \right]$$

$\gamma \in [0, 1)$: discount factor (value future reward less and less)



$$r_0 + \gamma \cdot r_1 + \gamma^2 \cdot r_2 + \dots + \gamma^h \cdot r_h + \dots$$

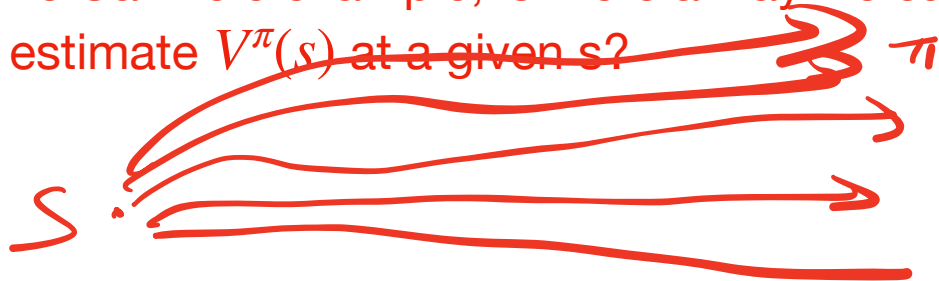
Performance of a policy π

Expected total reward of a policy π :

$$V^\pi(s) = \mathbb{E} \left[r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^h r_h + \dots \mid s_0 = s, \pi \right]$$

$\gamma \in [0,1)$: discount factor (value future reward less and less)

Q: think about the CartPole example, is there a way we can estimate $V^\pi(s)$ at a given s ?



Optimal policy

π^* : the policy that maximizes expected future reward at all states

$$V^*(s) \geq V^\pi(s), \forall s, \forall \pi$$

Optimal policy

π^* : the policy that maximizes expected future reward at all states

$$V^*(s) \geq V^\pi(s), \forall s, \forall \pi$$

Fact: π^* is deterministic, only depends on state (i.e., Markovian);
 π^* always exists for any infinite horizon discounted MDP

Optimal policy

π^* : the policy that maximizes expected future reward at all states

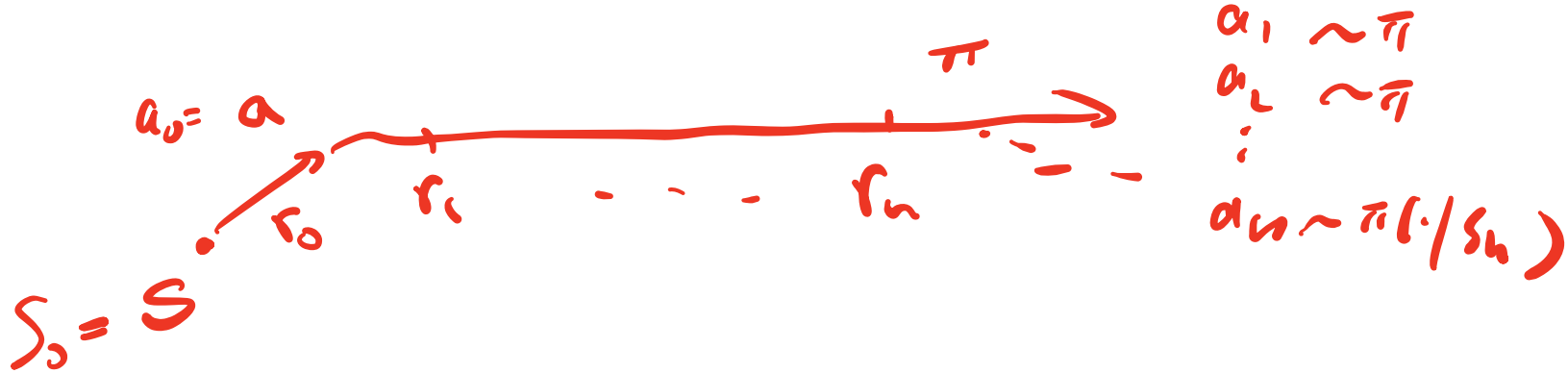
$$V^*(s) \geq V^\pi(s), \forall s, \forall \pi$$

Fact: π^* is deterministic, only depends on state (i.e., Markovian);
 π^* always exists for any infinite horizon discounted MDP

Q: what is the optimal policy when $\gamma = 0$?


State-action Q function

$$Q^\pi(s, a) = \mathbb{E} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^h r_h + \dots \mid s_0 = s, a_0 = a, \pi]$$



Outlines:

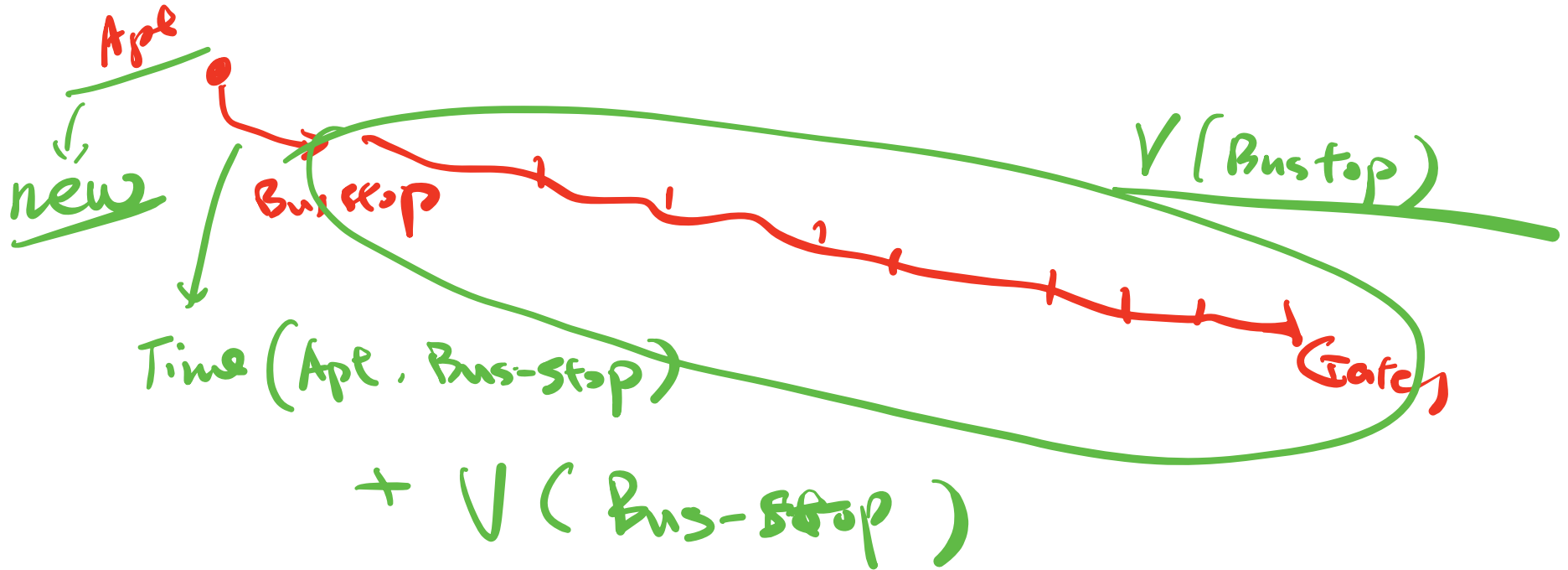
 1. Definitions of Markov Decision Process

 2. Value functions (V and Q functions)

3. Bellman equations

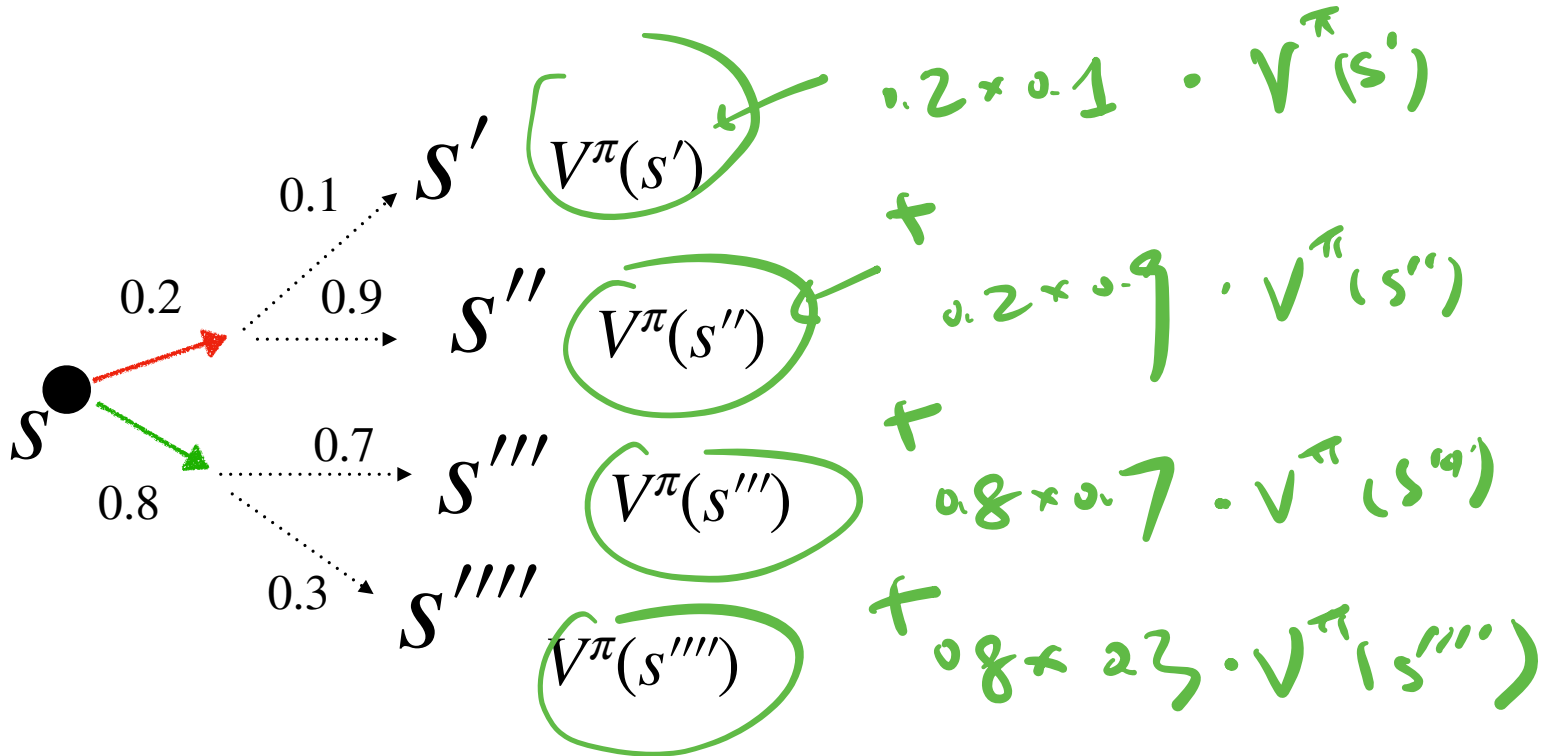
Can we quantify V / Q using one-step transition?

$$V^\pi(s) = \mathbb{E} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^h r_h + \dots | s_0 = s, a \sim \pi]$$



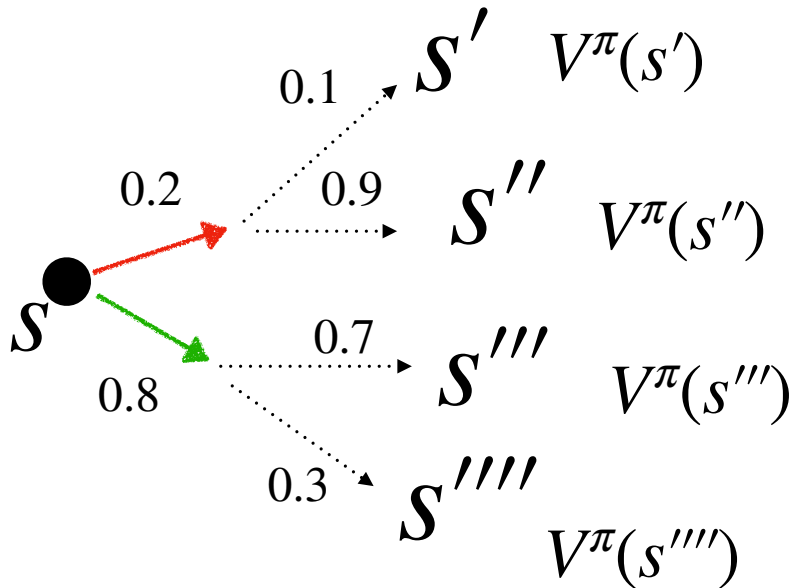
Can we quantify V / Q using one-step transition?

$$V^\pi(s) = \mathbb{E} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^h r_h + \dots | s_0 = s, a \sim \pi]$$



Can we quantify V / Q using one-step transition?

$$V^\pi(s) = \mathbb{E} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots + \gamma^h r_h + \dots | s_0 = s, a \sim \pi]$$



$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} [r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')]]$$

Bellman equation for value function

Can we quantify V / Q using one-step transition?

Your homework: understand the one-step relationship between V and Q

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^\pi(s')]$$

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot | s)} Q^\pi(s, a)$$

Summary:

- **Discounted infinite horizon MDP:**
 - State, action, policy, transition, reward (or cost), discount factor
 - **V function and Q function**
 - Key concept: **Bellman equation**