# Lecture 2: Markov Decision Processes

Kevin Chen and Zack Khan

Slides from David Silver

# Outline

1. Review of Last Lecture

2. Intro to MDPs

3. Markov Chains

4. Markov Reward Processes

5. Markov Decision Processes

6. Snippet of Bellman Expectation Equation for Markov Chain

## Review of Last Lecture

# Summary of Lecture 1

1. Reinforcement Learning (RL) is about an agent maximizing reward by interacting with its surrounding environment

2. RL has distinct advantages over other AI methods, but often requires more data or understanding of the problem

3. Agents take actions within an environment. Environment responds with rewards (or no reward) After an action, the agent moves into a new state of the environment



4. Figuring out how to tell an agent what actions to take, in order to maximize reward, is the key to reinforcement learning and creating a good AI

# Intro to MDPs

# Introduction to MDPs

- *Markov decision processes* formally describe an environment for reinforcement learning
- Where the environment is *fully observable*
- i.e. The current *state* completely characterises the process
- Almost all RL problems can be formalised as MDPs

# Markov Property

"The future is independent of the past given the present"

### Definition

A state $S_t$ is *Markov* if and only if

$$P[S_{t+1} \mid S_t] = P[S_{t+1} \mid S_1, ..., S_t]$$

- The state captures all relevant information from the history
- Once the state is known, the history may be thrown away
- i.e. The state is a sufficient statistic of the future

# State Transition Matrix

For a Markov state $s$ and successor state $s^j$, the *state transition probability* is defined by

$$\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$$

State transition matrix $\mathrm{P}$ defines transition probabilities from all states $s$ to all successor states $s^j$,

$$\mathcal{P} = \textit{from} \begin{bmatrix} \mathcal{P}_{11} & \ldots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \ldots & \mathcal{P}_{nn} \end{bmatrix}$$

where each row of the matrix sums to 1.
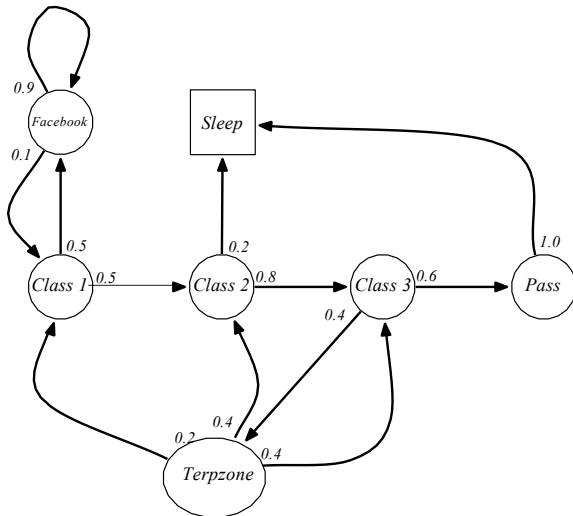
# Markov Chains

# Markov Chain

A Markov process is a memoryless random process, i.e. a sequence of random states $S_1$, $S_2$, ... with the Markov property.
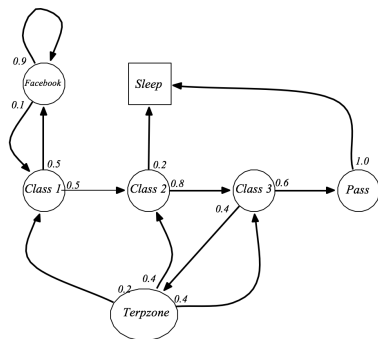
## Definition

A *Markov Process* (or *Markov Chain*) is a tuple $(S, P)$

- $S$ is a (finite) set of states
- $P$ is a state transition probability matrix,
  $P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$

# Example: Student Markov Chain
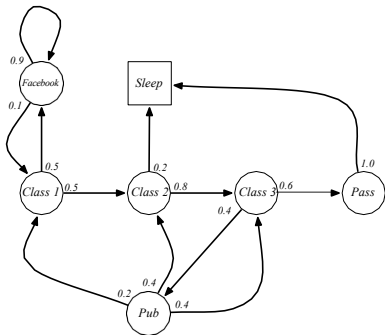
# Example: Student Markov Chain Episodes



Sample episodes for Student Markov Chain starting from $S_1 = $ C1

$$S_1, S_2, ..., S_T$$

- C1 C2 C3 Pass Sleep C1
- FB FB C1 C2 Sleep
- C1 C2 C3 Pub C2 C3 Pass Sleep
- C1 FB FB C1 C2 C3 Pub C1 FB FB

# Example: Student Markov Chain Transition Matrix



$$\mathcal{P} = \begin{array}{c|cccccc} & C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ \hline C1 & & 0.5 & & & & 0.5 & \\ C2 & & & 0.8 & & & & 0.2 \\ C3 & & & & 0.6 & 0.4 & & \\ Pass & & & & & & & 1.0 \\ Pub & 0.2 & 0.4 & 0.4 & & & & \\ FB & 0.1 & & & & & 0.9 & \\ Sleep & & & & & & & 1 \end{array}$$

## Markov Reward Processes
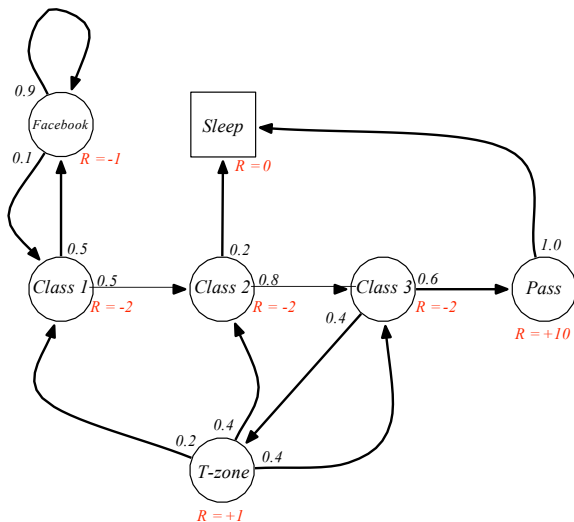
# Markov Reward Process

A Markov reward process is a Markov chain with values.

## Definition

A *Markov Reward Process* is a tuple $(S, P, R, \gamma)$

- $S$ is a finite set of states
- $P$ is a state transition probability matrix,
  $P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$
- $R$ is a reward function, $R_s = E[R_{t+1} \mid S_t = s]$
- $\gamma$ is a discount factor, $\gamma \in [0, 1]$

# Example: Student MRP

## Return

### Definition

The *return* $G_t$ is the total discounted reward from time-step $t$.

$$G_t = R_{t+1} + \gamma R_{t+2} + ... = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- The *discount* $\gamma \in [0, 1]$ is the present value of future rewards
- The value of receiving reward $R$ after $k + 1$ time-steps is $\gamma^k R$.
- This values immediate reward above delayed reward.
  - $\gamma$ close to 0 leads to short-term evaluation
  - $\gamma$ close to 1 leads to "far-sighted" evaluation

# Why discount?

Most Markov reward and decision processes are discounted. Why?

- Mathematically convenient to discount rewards
- Avoids infinite returns in cyclic Markov processes
- Uncertainty about the future may not be fully represented
- If the reward is financial, immediate rewards may earn more interest than delayed rewards
- Animal/human behaviour shows preference for immediate reward
- It is sometimes possible to use *undiscounted* Markov reward processes (i.e. $\gamma = 1$), e.g. if all sequences terminate.

# Value Function

The value function $v(s)$ gives the long-term value of state $s$

## Definition

The *state value function $v(s)$* of an MRP is the expected return starting from state $s$
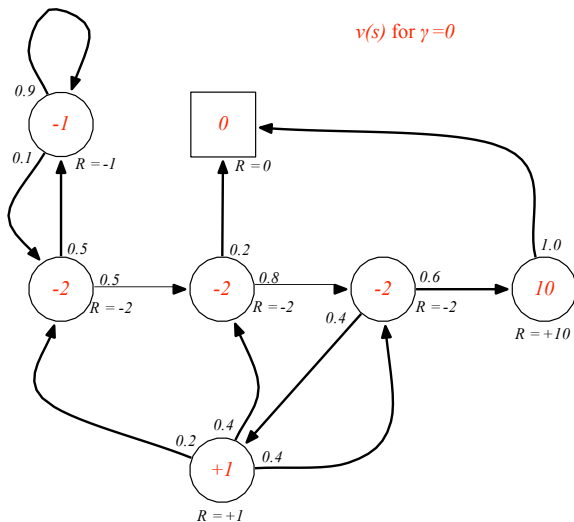
$$v(s) = \mathbb{E}[G_t \mid S_t = s]$$

# Example: Student MRP Returns

Sample returns for Student MRP:
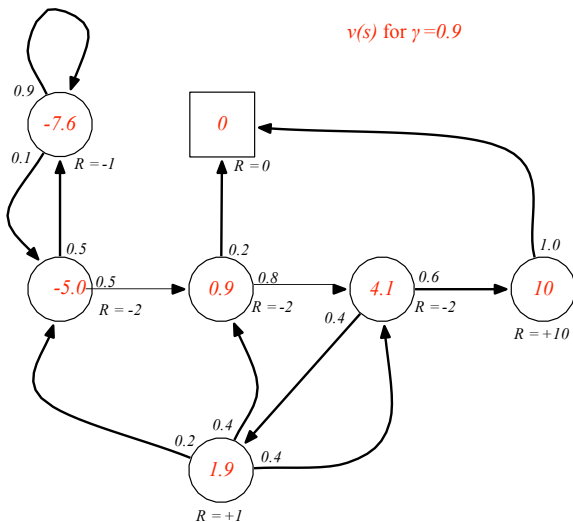Starting from $S_1 = $ C1 with $\gamma = 1/2$

$$G_1 = R_2 + \gamma R_3 + ... + \gamma^{T-2} R_T$$

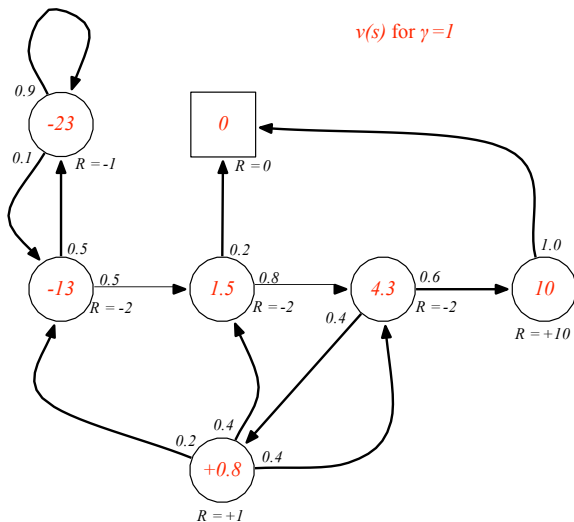| | | | |
|---|---|---|---|
| C1 C2 C3 Pass Sleep | $v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$ | $=$ | $-2.25$ |
| C1 FB FB C1 C2 Sleep | $v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$ | $=$ | $-3.125$ |
| C1 C2 C3 Pub C2 C3 Pass Sleep | $v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} \cdots$ | $=$ | $-3.41$ |
| C1 FB FB C1 C2 C3 Pub C1 ... | $v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} \cdots$ | $=$ | $-3.20$ |
| FB FB FB C1 C2 C3 Pub C2 Sleep | | | |

# Example: State-Value Function for Student MRP (1)



$v(s)$ for $\gamma = 0$

# Example: State-Value Function for Student MRP (2)



$v(s)$ for $\gamma=0.9$

# Example: State-Value Function for Student MRP (3)



$v(s)$ for $\gamma = 1$

# Markov Decision Processes

# Markov Decision Process

A Markov decision process (MDP) is a Markov reward process with decisions. It is an *environment* in which all states are Markov.
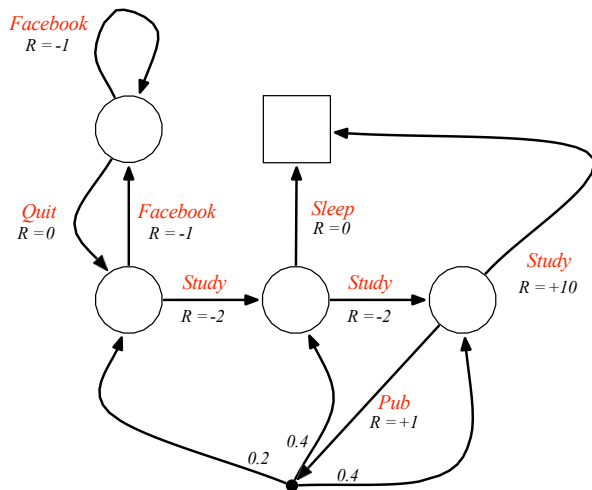
## Definition

A *Markov Decision Process* is a tuple $(S, A, P, R, \gamma)$

- $S$ is a finite set of states
- $A$ is a finite set of actions
- $P$ is a state transition probability matrix,
  $P_{ss'}^{a} = P[S_{t+1} = s' \mid S_t = s, A_t = a]$
- $R$ is a reward function, $R_s^a = E[R_{t+1} \mid S_t = s, A_t = a]$
- $\gamma$ is a discount factor $\gamma \in [0, 1]$.

# Example: Student MDP

# Policies (1)

## Definition

A *policy* $\pi$ is a distribution over actions given states,

$$\pi(a|s) = P[A_t = a \mid S_t = s]$$

- A policy fully defines the behaviour of an agent
- MDP policies depend on the current state (not the history)
- i.e. Policies are *stationary* (time-independent),
$$A_t \sim \pi(\cdot | S_t), \ \forall t > 0$$

# Policies (2)

- Given an MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$ and a policy $\pi$
- The state sequence $S_1, S_2, \ldots$ is a Markov process $(\mathcal{S}, \mathcal{P}^\pi)$
- The state and reward sequence $S_1, R_2, S_2, \ldots$ is a Markov reward process $(\mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma)$
- where

$$\mathcal{P}^\pi_{s,s'} = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{P}^a_{ss'}$$

$$\mathcal{R}^\pi_s = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{R}^a_s$$

# Value Function

### Definition

The *state-value function* $v_\pi(s)$ of an MDP is the expected return starting from state $s$, and then following policy $\pi$
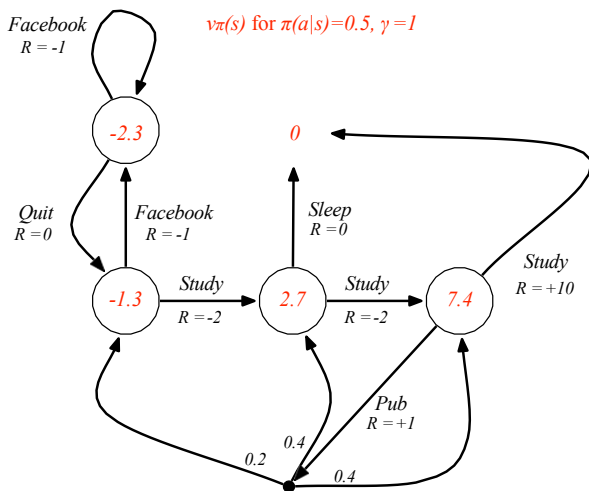
$$v_\pi(s) = E_\pi[G_t \mid S_t = s]$$

### Definition

The *action-value function* $q_\pi(s, a)$ is the expected return starting from state $s$, taking action $a$, and then following policy $\pi$

$$q_\pi(s, a) = E_\pi[G_t \mid S_t = s, A_t = a]$$

# Example: State-Value Function for Student MDP



$v_\pi(s)$ for $\pi(a|s)=0.5$, $\gamma =1$

## Optimal Value Function

### Definition

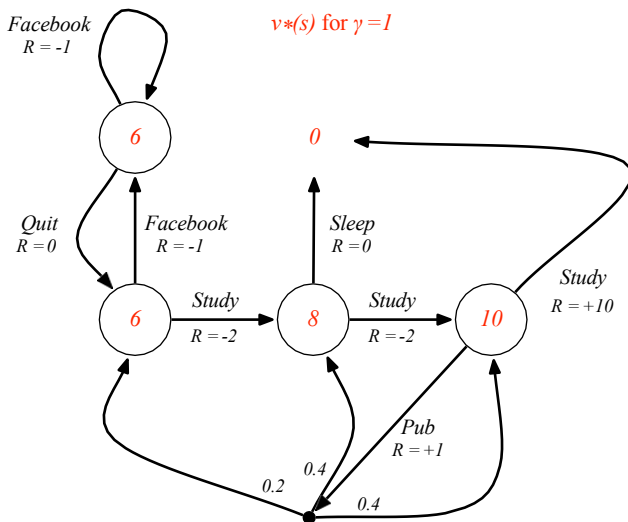The *optimal state-value function* $v_*(s)$ is the maximum value function over all policies

$$v_*(s) = \max_\pi v_\pi(s)$$

The *optimal action-value function* $q_*(s, a)$ is the maximum action-value function over all policies
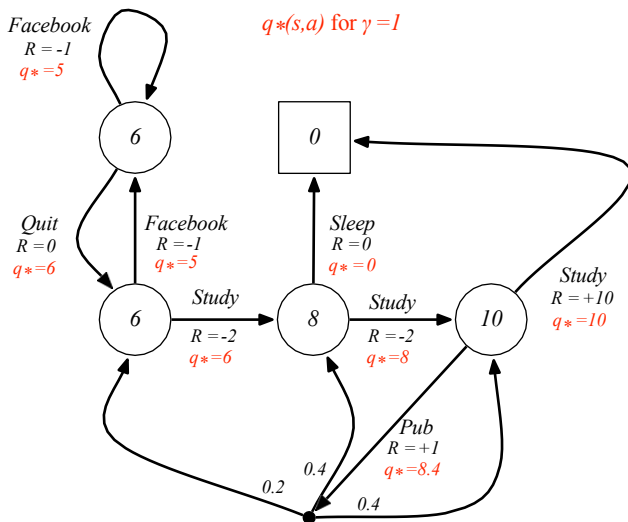
$$q_*(s, a) = \max_\pi q_\pi(s, a)$$

- The optimal value function specifies the best possible performance in the MDP.
- An MDP is "solved" when we know the optimal value fn.

# Example: Optimal Value Function for Student MDP



$v*(s)$ for $\gamma = 1$

Facebook
R = -1

6

0

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

6    Study    8    Study    10
     R = -2         R = -2

Pub
R = +1

0.4

0.2

0.4

# Example: Optimal Action-Value Function for Student MDP

## Optimal Policy

Define a partial ordering over policies

$$\pi \geq \pi \text{ if } v_\pi(s) \geq v_{\pi^i}(s), \ \forall s$$

### Theorem

*For any Markov Decision Process*

- *There exists an optimal policy $\pi_*$ that is better than or equal to all other policies, $\pi_* \geq \pi, \ \forall \pi$*
- *All optimal policies achieve the optimal value function, $v_{\pi_*}(s) = v_*(s)$*
- *All optimal policies achieve the optimal action-value function, $q_{\pi_*}(s, a) = q_*(s, a)$*
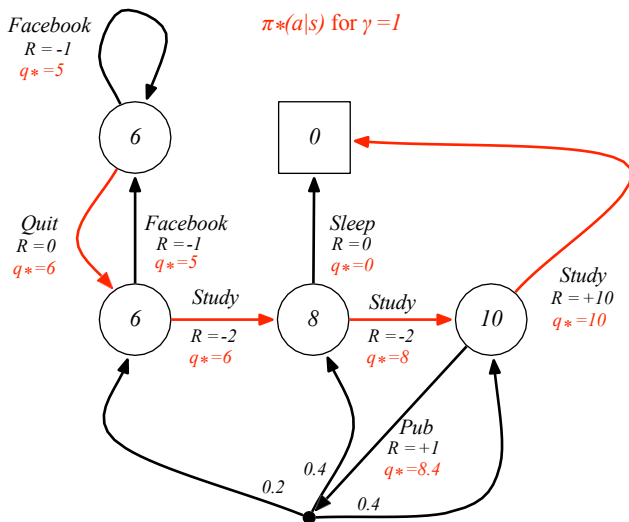
# Finding an Optimal Policy

An optimal policy can be found by maximising over $q_*(s, a)$,

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in \mathcal{A}}{\text{argmax}} \; q_*(s, a) \\ 0 & otherwise \end{cases}$$

- There is always a deterministic optimal policy for any MDP
- If we know $q_*(s, a)$, we immediately have the optimal policy

# Example: Optimal Policy for Student MDP

Bellman Expectation Equation for Markov Chain

## Bellman Equation for MRPs

The value function can be decomposed into two parts:

- immediate reward $R_{t+1}$
- discounted value of successor state $\gamma v(S_{t+1})$

$$
\begin{aligned}
v(s) &= \mathbb{E}\left[G_t \mid S_t = s\right] \\
&= \mathbb{E}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s\right] \\
&= \mathbb{E}\left[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s\right] \\
&= \mathbb{E}\left[R_{t+1} + \gamma G_{t+1} \mid S_t = s\right] \\
&= \mathbb{E}\left[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s\right]
\end{aligned}
$$

Questions?