

Lecture 8: Markov Decision Process

Ziyu Shao

School of Information Science and Technology
ShanghaiTech University

April 27 & 29, 2020

Outline

1 Introduction

2 Markov Reward Process

3 Markov Decision Process

4 References

Outline

1 Introduction

2 Markov Reward Process

3 Markov Decision Process

4 References

Markov Decision Process

- Markov decision processes formally describe an environment for reinforcement learning
- Where the environment is fully observable
- i.e. The current state completely characterizes the process
- Almost all RL problems can be formalized as MDPs, e.g
 - ▶ Optimal control primarily deals with continuous MDPs
 - ▶ Partially observable problems can be converted into MDPs
 - ▶ Bandits are MDPs with one state

POMDP

Markov Property

Definition

A state S_t is Markovian if and only if

$$\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t]$$

- “The future is independent of the past given the present”
- 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

Markov Chain

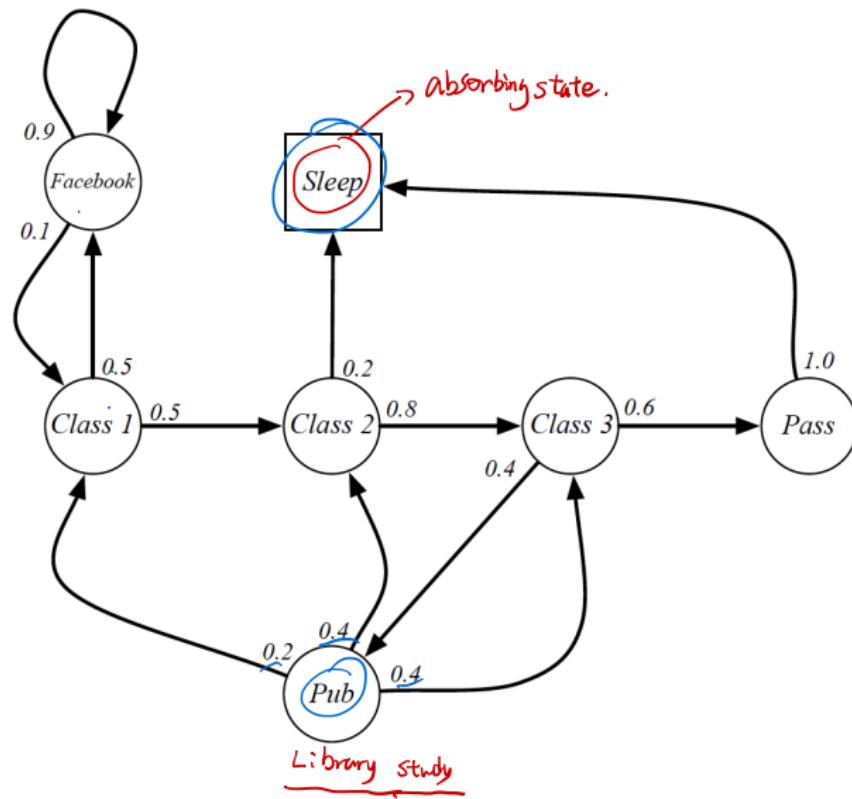
Definition

A discrete-time Markov chain is a tuple $\langle \mathcal{S}, \mathcal{P} \rangle$

- \mathcal{S} is a (finite) set of states
- \mathcal{P} is a state transition probability matrix

$$\mathcal{P}_{s,s'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$$

Example: Student Markov Chain

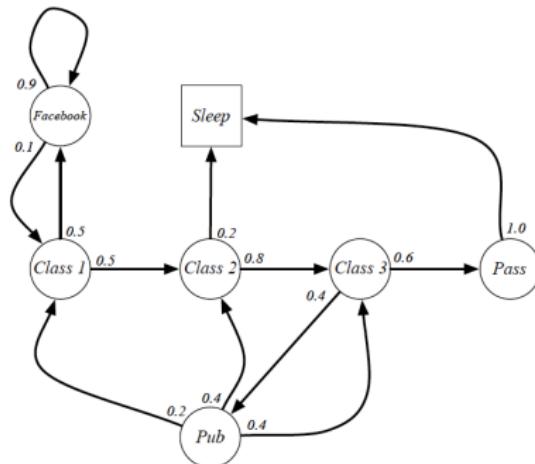


Example: Student Markov Chain Episodes

Sample path

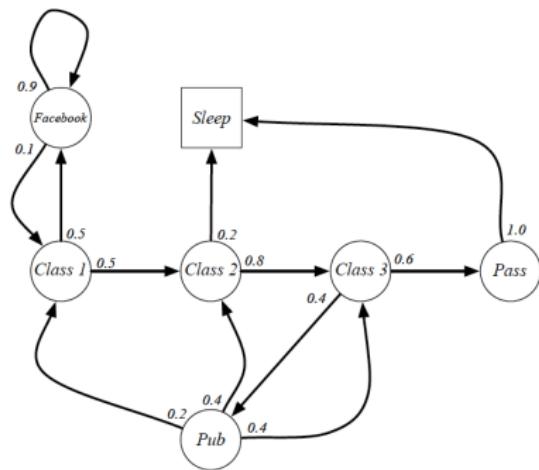
Sample episodes for Student Markov Chain starting from $S_1 = C_1$

S_1, S_2, \dots, 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 FB C1 C2 C3 Pub C2 Sleep

Example: Student Markov Chain Transition Matrix



$$\mathcal{P} = \begin{bmatrix} C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ C1 & 0.5 & 0.8 & 0.6 & 0.4 & 0.5 & 0.2 \\ C2 & 0.2 & 0.4 & 0.6 & 0.4 & 0.4 & 1.0 \\ C3 & 0.1 & 0.4 & 0.4 & 0.9 & 0.1 & 0 \\ Pass & & & & & & \\ Pub & & & & & & \\ FB & & & & & & \\ Sleep & & & & & & \end{bmatrix}$$

Sum of each Row =
Stochastic matrix.

Outline

1 Introduction

2 Markov Reward Process

3 Markov Decision Process

4 References

Markov Reward Process

A Markov reward process is a Markov chain with values.

Definition

A Markov Reward Process is a tuple $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- \mathcal{S} is a (finite) set of states
- \mathcal{P} is a state transition probability matrix

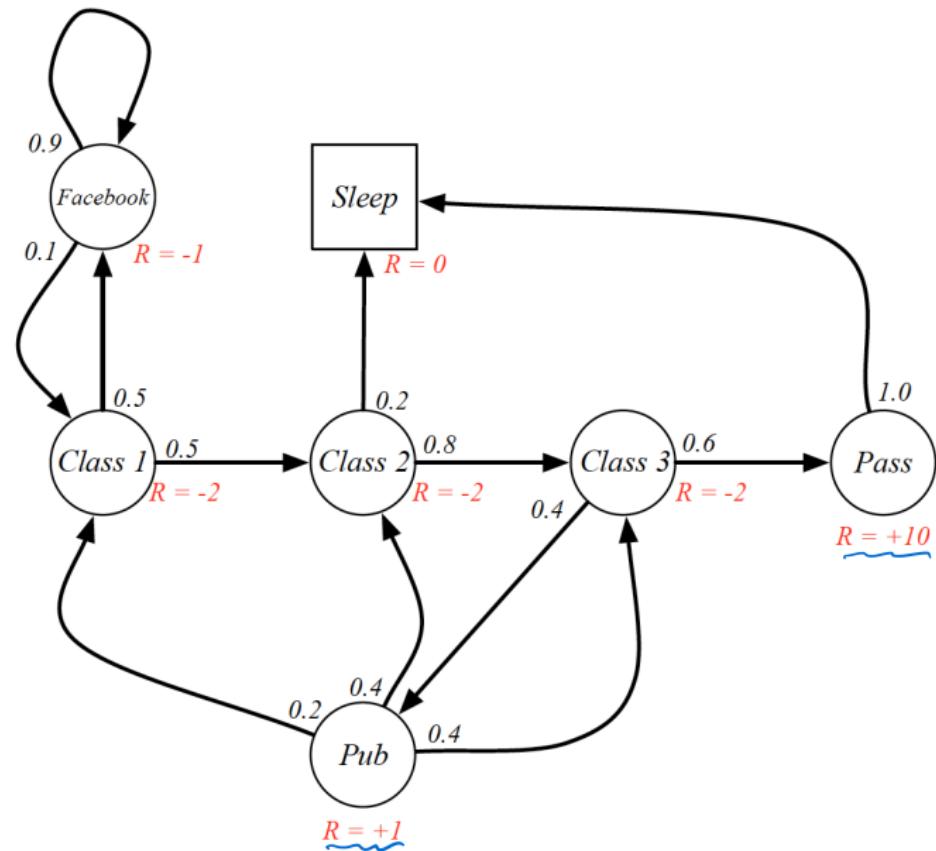
$$\mathcal{P}_{s,s'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$$

- \mathcal{R} is a reward function,

$$\mathcal{R}_s = \mathbb{E}[R_{t+1} | S_t = s]$$

- γ 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 .

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

- 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.
 - γ close to 0 leads to "myopic" evaluation
 - γ 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 behavior 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

$$v(s) = \mathbb{E}[G_t | S_t = s]$$

Example: Student MRP Returns

Sample **returns** for Student MRP:

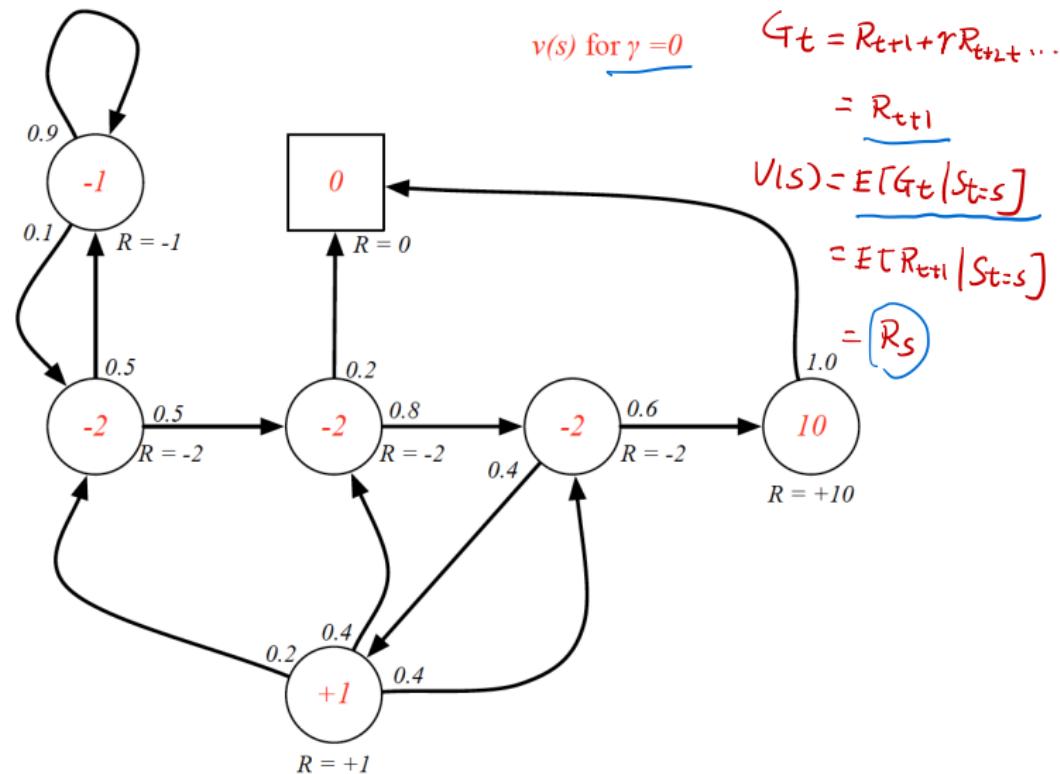
Starting from $S_1 = C1$ with $\gamma = \frac{1}{2}$

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

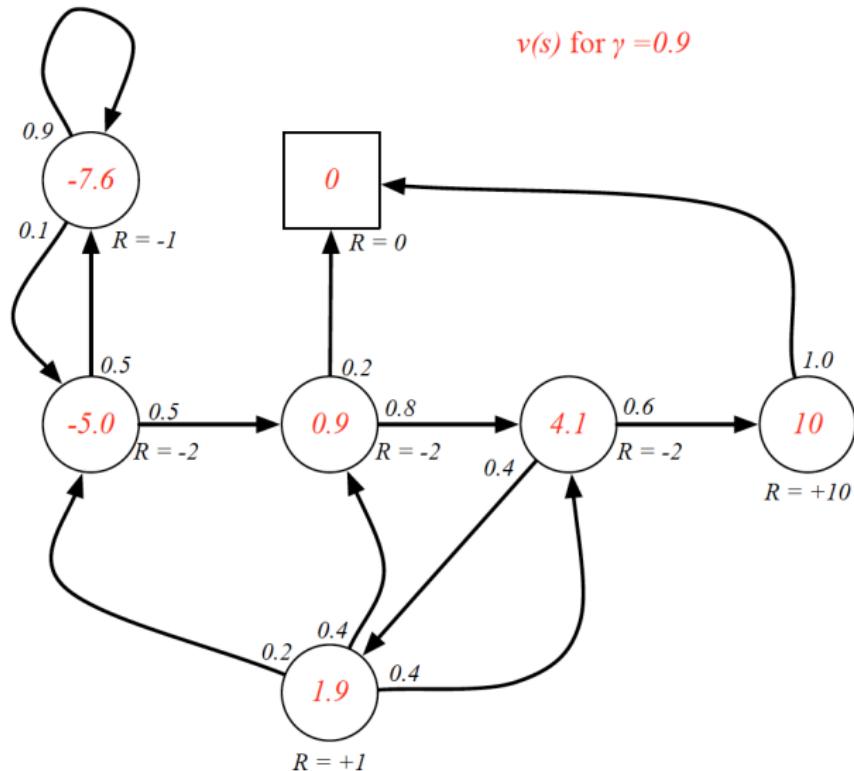
S_1	S_2	S_3	S_4	S_5	S_6	\dots
C_1	C1	C2	C3	Pass	Sleep	
C_1	C1	FB	FB	C1	C2	Sleep
C_1	C1	C2	C3	Pub	C2	C3
C_1	C1	FB	FB	C1	C2	C3
G	FB	FB	FB	C1	C2	C3

$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$	=	-2.25
$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$	=	-3.125
$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.41
$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.20
\dots		

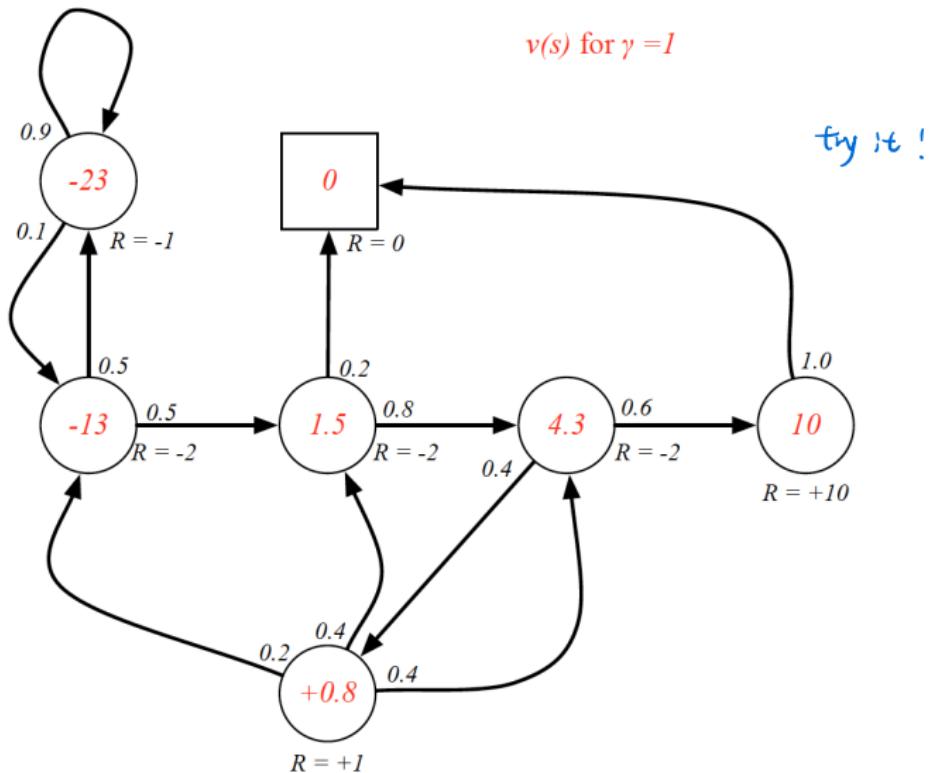
Example: State-Value Function for Student MRP



Example: State-Value Function for Student MRP



Example: State-Value Function for Student MRP



Bellman Equation for MRPs

The value function can be decomposed into two parts:

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

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) = R_{t+1} + \gamma G_{t+1}$$

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

$$V(S_t) = \mathbb{E}[R_{t+1} + \gamma V(S_{t+1}) | S_t]$$

Bellman Equation for MRPs

Adam's theorem + Markov.

$$1^{\circ}. \underline{V(S) = E[G_t | S_t=s]} , \underline{V(S_t) = E[G_t | S_t]} , \underline{V(S_{t+1}) = E[G_{t+1} | S_{t+1}]}$$

$$2^{\circ}. \text{ Adam's Law. } \underline{E[E(Y|X)] = E(Y)}$$

$$\text{Adam's Law with extra conditioning. } \hat{E}(\cdot) = E(\cdot | Z) \Rightarrow \hat{E}[\hat{E}(Y|X)] = \hat{E}(Y)$$

$$\Leftrightarrow E[E(Y|X, Z)|Z] = E(Y|Z)$$

Let $Y = G_{t+1}$, $X = S_{t+1}$, $Z = S_t$;
Then we have

$$\underline{E[E(G_{t+1} | S_{t+1}, S_t) | S_t] = E(G_{t+1} | S_t)} .$$

Markov Property

$$\underline{- E[E[G_{t+1} | S_{t+1}] | S_t]} = \underline{E[V(S_{t+1}) | S_t]}$$

$$\text{thus } \underline{E[G_{t+1} | S_t] = E[V(S_{t+1}) | S_t]}$$

$$\Rightarrow \underline{E[G_{t+1} | S_t=s] = E[V(S_{t+1}) | S_t=s]} .$$

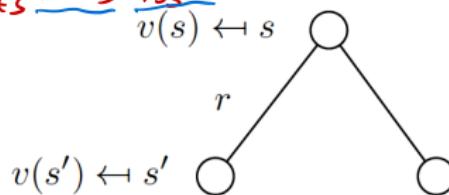
$$3^{\circ}. \underline{V(S) = E[G_t | S_t=s]} = \underline{E[R_{t+1} + \gamma G_{t+1} | S_t=s]} = \underline{E[R_{t+1} | S_t=s]} + \underline{\gamma E[G_{t+1} | S_t=s]} \\ = \underline{E[R_{t+1} | S_t=s]} + \underline{\gamma E[V(S_{t+1}) | S_t=s]} \\ = \underline{E[R_{t+1} + \gamma V(S_{t+1}) | S_t=s]}$$

Q.E.D.

Bellman Equation for MRPs

Bellman Equation for MRPs

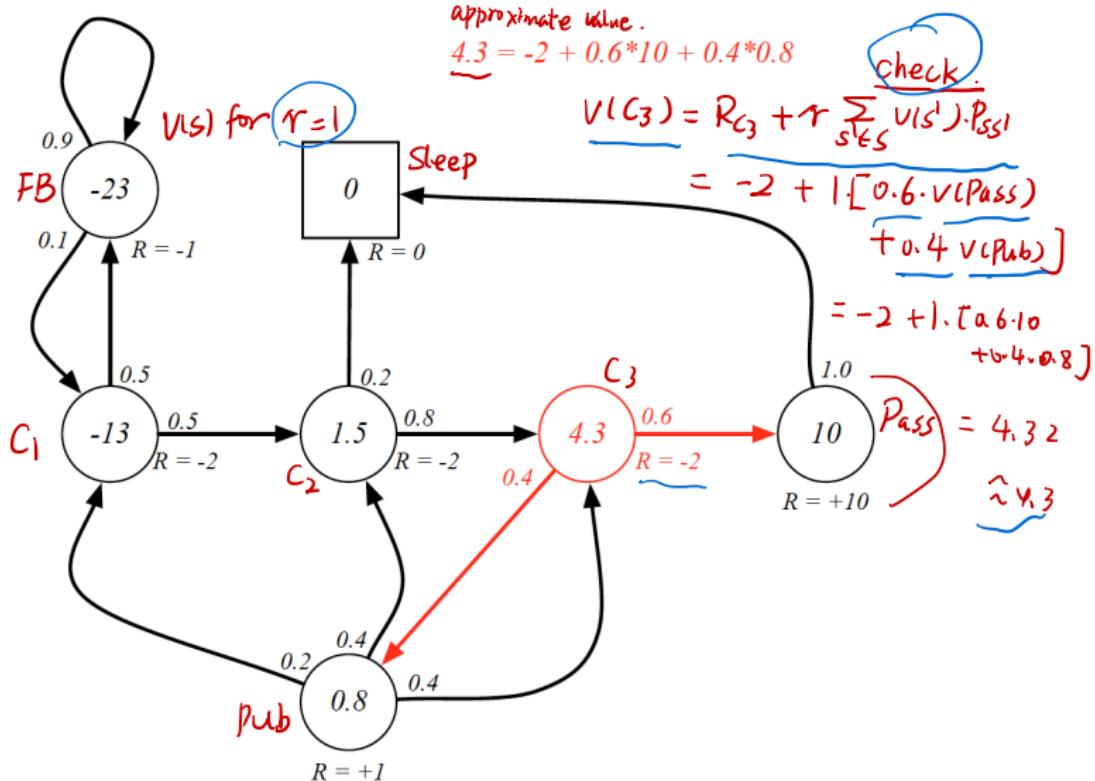
$$\begin{aligned} v(s) &= \underbrace{E[R_{t+1} | S_{t+1}=s]}_{R_s} + \gamma \underbrace{E[V(S_{t+1}) | S_{t+1}=s]}_{\text{Lote.}} \\ &= R_s + \gamma \sum_{s' \in S} E[V(S_{t+1}) | S_t=s, S_{t+1}=s'] p(S_{t+1}=s' | S_t=s) \\ v(s) &= \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t=s] \\ &= R_s + \gamma \sum_{s' \in S} v(s') \cdot p_{ss'} \end{aligned}$$



$$v(s) = R_s + \gamma \sum_{s' \in S} p_{ss'} v(s')$$

$$V(S_t) = E[R_{t+1} + \gamma V(S_{t+1}) | S_t]$$

Example: Bellman Equation for Student MRP



Bellman Equation in Matrix Form

The Bellman equation can be expressed concisely using matrices,

$$v = \underbrace{\mathcal{R} + \gamma \mathcal{P} v}$$

where v is a column vector with one entry per state

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} \mathcal{R}_1 \\ \vdots \\ \mathcal{R}_n \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

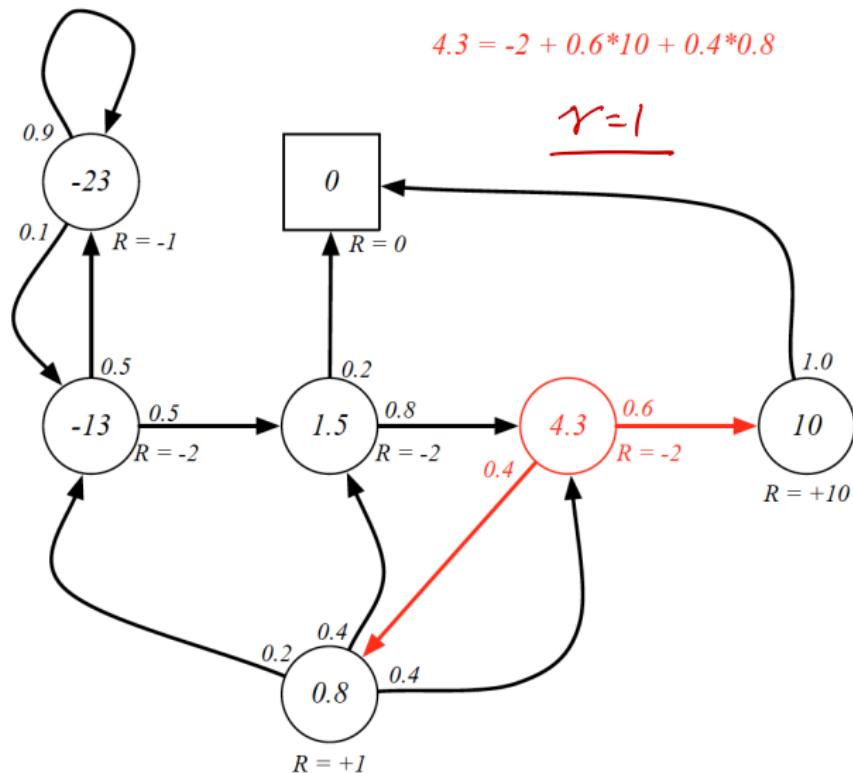
Solving the Bellman Equation

- The Bellman equation is a linear equation
- It can be solved directly:

$$\begin{aligned}v &= \mathcal{R} + \gamma \mathcal{P}v \\(I - \gamma \mathcal{P})v &= \mathcal{R} \\v &= (I - \gamma \mathcal{P})^{-1} \mathcal{R}\end{aligned}$$

- Computational complexity is $O(n^3)$ for n states
- Direct solution only possible for small MRPs
- There are many iterative methods for large MRPs, e.g.
 - Dynamic programming
 - Monte-Carlo evaluation
 - Temporal-Difference learning

Example: Matrix Solution for Student MRP



Example: Matrix Solution for Student MRP

1^o. state transition probability matrix P (Refer to page 9)

2^o. $R = \begin{bmatrix} -2 \\ -2 \\ -2 \\ 10 \\ 1 \\ -1 \\ 0 \end{bmatrix}$ $\underbrace{\begin{array}{c} c_1 \\ c_2 \\ c_3 \\ \text{Pass} \\ \text{Rub} \\ \text{FB} \\ \text{Sleep} \end{array}}$

3^o. $V = (I - rP)^{-1}R$ $\underbrace{\phantom{V = (I - rP)^{-1}R}_{(e-6)}}$

$$r = 1.5 \quad V = \begin{bmatrix} -12.529 \\ 1.4568 \\ 4.3210 \\ 10.0000 \\ 0.8025 \\ -22.5427 \\ 0 \end{bmatrix}$$

V = (I - rP)^{-1}R

$$r = 0.9 \quad V = \begin{bmatrix} -5.0127 \\ 0.9427 \\ 4.0870 \\ 10.0000 \\ 1.9084 \\ -7.6376 \\ 0 \end{bmatrix}$$

V = (I - rP)^{-1}R

Example: Matrix Solution for Student MRP

Outline

1 Introduction

2 Markov Reward Process

Markov Chain + ^{State.}
Reward



3 Markov Decision Process

-- - - - . - + ^{action}

4 References

Markov Decision Process

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

Definition

A Markov Decision Process is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- \mathcal{S} is a (finite) set of states
- \mathcal{A} is a finite set of actions
- \mathcal{P} is a state transition probability matrix

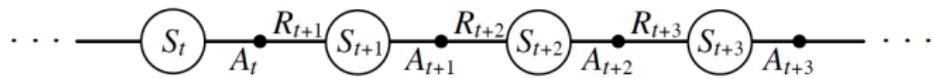
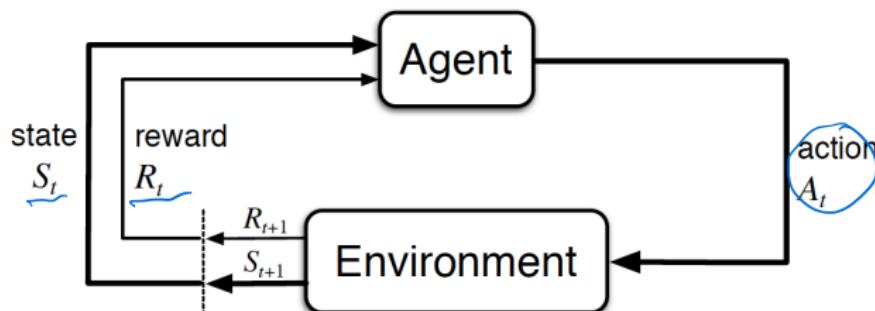
$$\mathcal{P}_{s,s'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

- \mathcal{R} is a reward function,

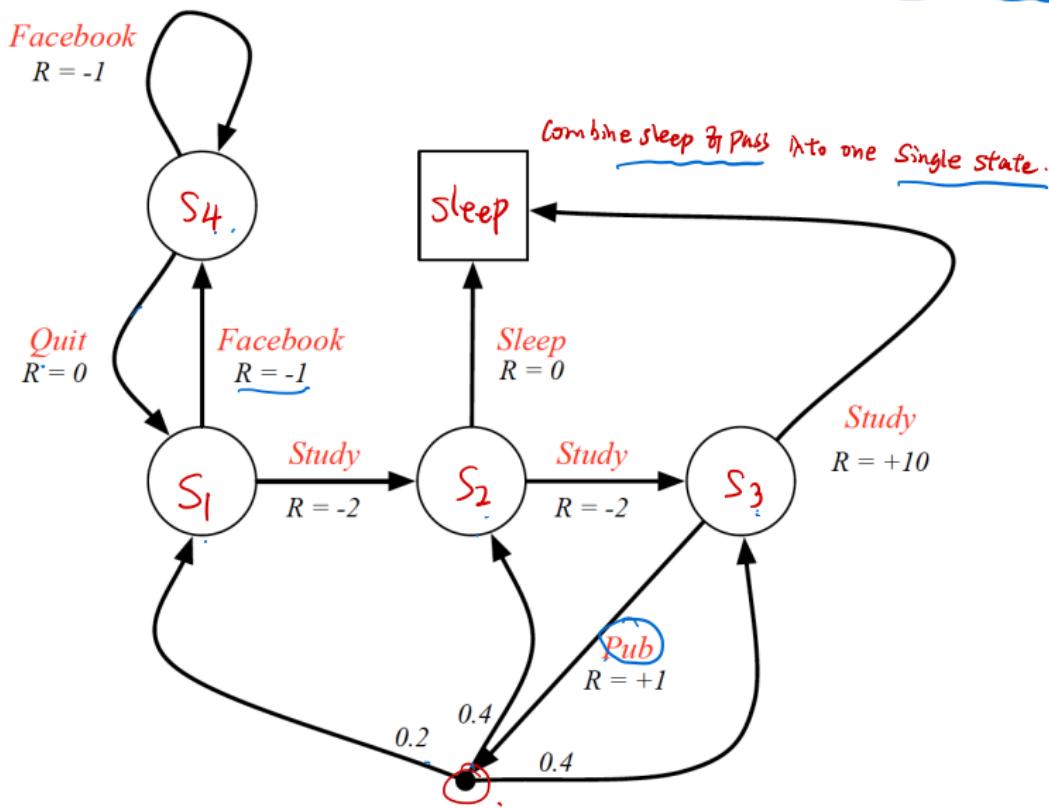
$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

- γ is a discount factor, $\gamma \in [0, 1]$

The Agent-Environment Interface



Example: Student MDP



Policy

Definition

A *policy* π is a distribution over actions given states,

$$\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$$
$$\sum_{a \in A} \pi(a|s) = 1$$

- 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),
 $A_t \sim \pi(\cdot|S_t), \forall t > 0$

Policy

$$1^o. P_{S_1, S'}^\pi = \underbrace{P[S_{t+1} = s' | S_t = s]}_{\text{LOT P}} = \sum_{a \in A} P[S_{t+1} = s' | A_t = a, S_t = s] p(A_t = a | S_t = s)$$

$$= \sum_{a \in A} P_{ss'}^a \cdot \pi(a|s)$$

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

$$1^o. \mathcal{P}_{s, s'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^a$$

$$2^o. \mathcal{R}_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a$$

$$R_s^\pi = E_\pi [R_{t+1} | S_t = s] \stackrel{\text{LOT E}}{=} \sum_{a \in A} E[R_{t+1} | S_t = s, A_t = a] p(A_t = a | S_t = s)$$

$$= \sum_{a \in A} R_s^a \cdot \pi(a|s)$$

Value Function

Adam's Law + Markov Property

Definition

The state-value function $v_\pi(s)$ of an MDP is the expected return starting from state s , and then following policy π

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s]$$

$$v_\pi(S_t) = \mathbb{E}_\pi [G_t \mid S_t]$$

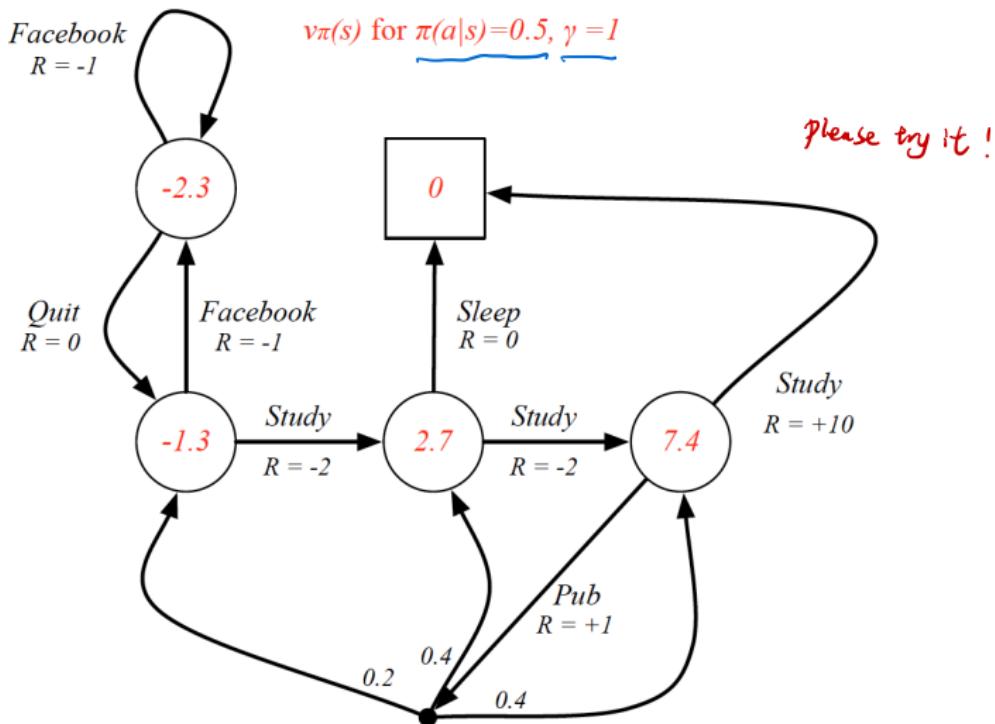
Definition

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

$$q_\pi(s, a) = \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a]$$

$$q_\pi(S_t, A_t) = \mathbb{E}_\pi [G_t \mid S_t, A_t]$$

Example: State-Value Function for Student MDP



Bellman Expectation Equation

The state-value function can again be decomposed into immediate reward plus discounted value of successor state,

$$V_{\pi}(S_t) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t]$$
$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

same proof as

MRP

You can try it!

The action-value function can similarly be decomposed,

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

proof next

$$q_{\pi}(S_t, A_t) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t, A_t]$$

Bellman Expectation Equation

1^o. $q_{\pi}(s, a) = E_{\pi}[G_t | S_t=s, A_t=a]$

$q_{\pi}(S_t, A_t) = E_{\pi}[G_t | S_t, A_t]$

$q_{\pi}(S_{t+1}, A_{t+1}) = E_{\pi}[G_{t+1} | S_{t+1}, A_{t+1}]$

2^o. Adam's Law with extra conditioning

$E[Y|Z] = E[E[Y|X, Z]|Z]$

Let $Y = G_{t+1}$, $Z = (S_t, A_t)$, $X = (S_{t+1}, A_{t+1})$,

then we have

$E[G_{t+1} | S_t, A_t] = E[E[E[G_{t+1} | S_{t+1}, A_{t+1}, S_t, A_t] | S_t, A_t]$

Markov Property

$= E[E[G_{t+1} | S_{t+1}, A_{t+1}] | S_t, A_t]$

$= E[q_{\pi}(S_{t+1}, A_{t+1}) | S_t, A_t]$

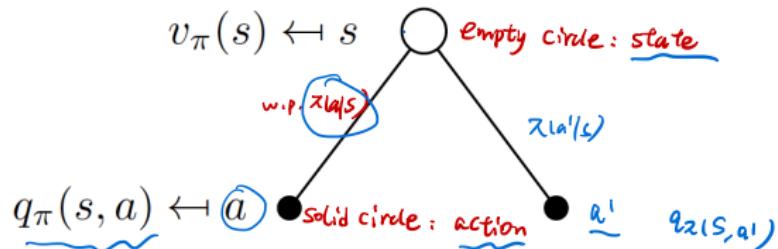
$\Rightarrow E_{\pi}[G_{t+1} | S_t=s, A_t=a] = E_{\pi}[q_{\pi}(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]$

Bellman Expectation Equation

3°. Then we have

$$\begin{aligned} q_{\pi}(s, a) &= \underline{E_{\pi}[G_t | S_t=s, A_t=a]} \\ G_t &= R_{t+1} + rG_{t+1} \\ &= \underline{E_{\pi}[R_{t+1} + rG_{t+1} | S_t=s, A_t=a]} \\ &= \underline{E_{\pi}[R_{t+1} | S_t=s, A_t=a]} + r \underline{E_{\pi}[G_{t+1} | S_t=s, A_t=a]} \\ &= \underline{E_{\pi}[R_{t+1} | S_t=s, A_t=a]} + r \underline{E_{\pi}[q_{\pi}(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]} \\ &= \underline{E_{\pi}[R_{t+1} + r q_{\pi}(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]} \end{aligned}$$

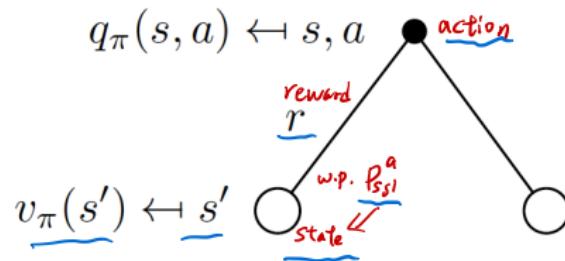
Bellman Expectation Equation for V^π



$$\underline{v_\pi(s)} = \sum_{a \in A} \pi(a|s) q_\pi(s, a)$$

$$V_\pi(s) = \underline{E_\pi[G_t | S_t=s]} \stackrel{\text{LoTE}}{=} \sum_{a \in A} \underline{E_\pi[G_t | S_t=s, A_t=a]} \underline{P(A_t=a | S_t=s)}$$
$$= \sum_{a \in A} \underline{q_\pi(s, a)} \underline{\pi(a|s)}$$

Bellman Expectation Equation for Q^π



$$q_\pi(s, a) = \underbrace{\mathcal{R}_s^a}_{\text{---}} + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \underbrace{v_\pi(s')}_{\text{---}}$$

Proof next

Bellman Expectation Equation for Q^π

$$1^o. E[\underbrace{q_\pi(s_{t+1}, a_{t+1})}_{\text{Markov Property}} | \underbrace{s_{t+1}=s', s_t=s, a_t=a}]$$

Markov Property

$$= E[\underbrace{q_\pi(s_{t+1}, a_{t+1})}_{\text{Markov Property}} | s_{t+1}=s']$$

LOTE

$$= \sum_{a \in A} E[q_\pi(s_{t+1}, a_{t+1}) | s_{t+1}=s', a_{t+1}=a] \cdot P(a_{t+1}=a | s_{t+1}=s')$$

$$= \sum_{a \in A} \underbrace{q_\pi(s', a)}_{\text{LOTE}} \underbrace{\pi(a | s')}_{\text{LOTE}}$$

$$= \underbrace{v_\pi(s')}$$

$$2^o. \text{ Then we have } E[q_\pi(s_{t+1}, a_{t+1}) | s_t=s, a_t=a]$$

$$\begin{aligned} &= \sum_{s' \in S} E[q_\pi(s_{t+1}, a_{t+1}) | s_{t+1}=s', s_t=s, a_t=a] \cdot P(s_{t+1}=s' | s_t=s, a_t=a) \\ &= \sum_{s' \in S} \underbrace{v_\pi(s')}_{\text{LOTE}} \cdot \underbrace{P_{ss'}^a}_{\text{LOTE}} \end{aligned}$$

Bellman Expectation Equation for Q^π

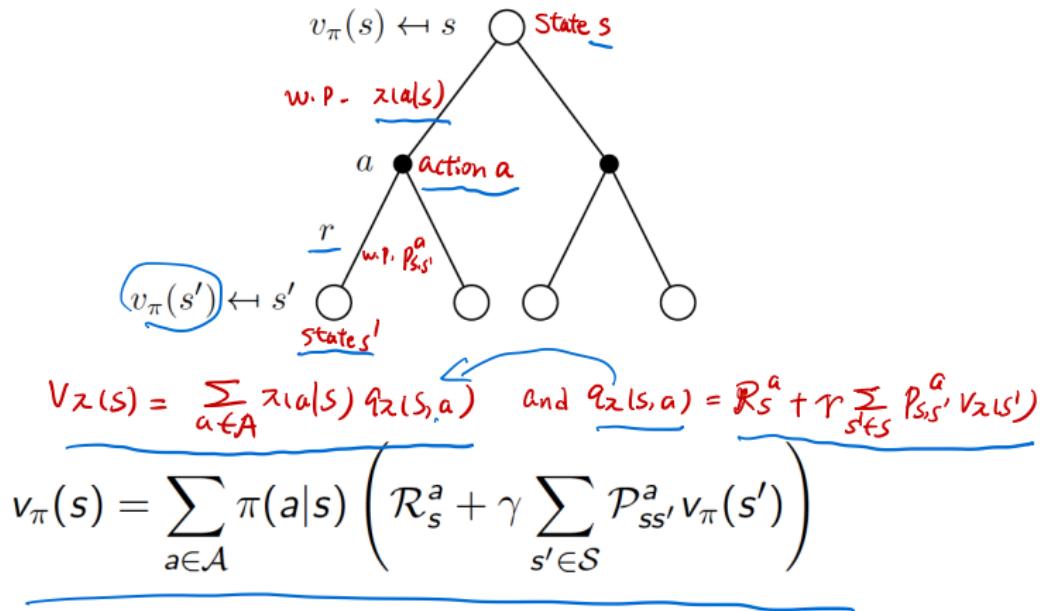
3^o. Thus $q_\pi(s, a) = \underbrace{E_\pi [R_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1}) | s_t=s, a_t=a]}$

$$= \underbrace{E[R_{t+1} | s_t=s, a_t=a]} + \gamma \underbrace{E[q_\pi(s_{t+1}, a_{t+1}) | s_t=s, a_t=a]}$$

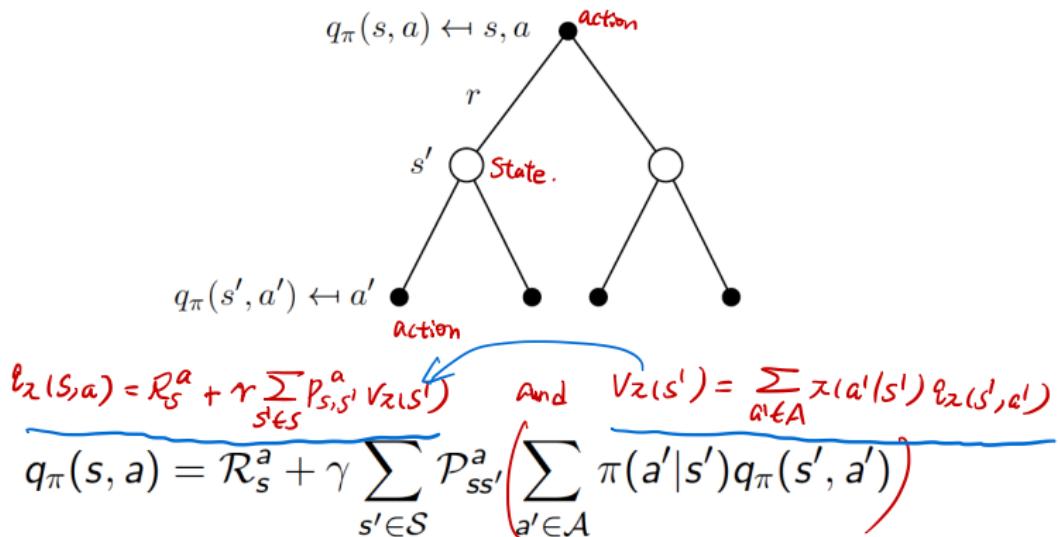
2^o.

$$= \underbrace{r_s^a}_{\text{1}} + \gamma \cdot \underbrace{\sum_{s' \in S} p_{s,s'}^a \cdot r_{\pi}(s')}$$

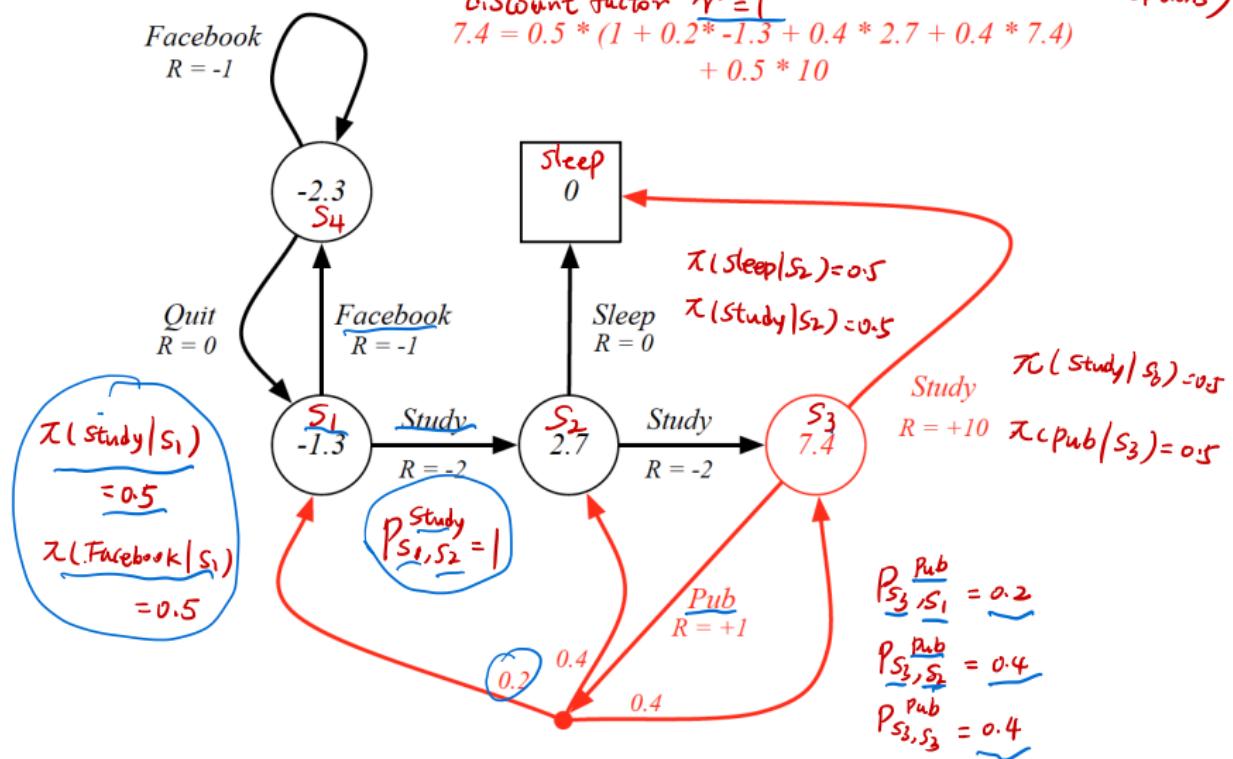
Bellman Expectation Equation for v_π



Bellman Expectation Equation for q_{π}



Example: Bellman Expectation Equation in Student MDP



Example: Bellman Expectation Equation in Student MDP

$$1^{\text{st}}. \quad V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left[R_s^a + \gamma \sum_{s' \in S} p_{s,s'}^a V_{\pi}(s') \right]$$

$$U_j \triangleq V_{\pi}(s_j)$$

$$2^{\text{nd}}. \quad \underbrace{V_1}_{\approx V_{\pi}(s_1)} = \underbrace{0.5 \left(\underbrace{R_{s_1}^{\text{Study}}}_{\uparrow \pi(\text{Study}|s_1)} + \underbrace{1.1 V_{\pi}(s_2)}_{\uparrow \pi(\text{Facebook}|s_1)} \right)}_{= 0.5 (-2 + U_2)} + \underbrace{0.5 \left(\underbrace{R_{s_1}^{\text{Facebook}}}_{\uparrow \pi(\text{Facebook}|s_1)} + \underbrace{1.1 V_{\pi}(s_4)}_{\uparrow \pi(\text{Facebook}|s_1)} \right)}_{= 0.5 (-1 + U_4)}$$

$$\underline{U_2 = V_{\pi}(s_2)} = 0.5 (-2 + U_3) + 0.5 (0 + 0)$$

$$\underline{U_3 = V_{\pi}(s_3)} = 0.5 (1 + 0.2U_1 + 0.4U_2 + 0.4U_3) + 0.5 (0 + 0)$$

$$\underline{U_4 = V_{\pi}(s_4)} = 0.5 (0 + U_1) + 0.5 (-1 + U_4)$$

Example: Bellman Expectation Equation in Student MDP

3^o. then we have

$$\begin{aligned}V_1 &= -1.3 \\V_2 &= 2.7 \\V_3 &= 7.4 \\V_4 &= -2.3\end{aligned}$$

4^o. We have obtained the state value. Now we turn to action. State-action value.

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{s,s'} V_{\pi}(s')$$

$$\Rightarrow q_{\pi}(s_1, \text{study}) = -2 + 1 \cdot V_2 = -2 + 2.7 = 0.7$$

$$q_{\pi}(s_1, \text{Facebook}) = -1 + 1 \cdot V_4 = -1 - 2.3 = -3.3$$

$$q_{\pi}(s_2, \text{sleep}) = 0 + 0 = 0.$$

$$q_{\pi}(s_2, \text{Study}) = -2 + 1 \cdot V_3 = -2 + 7.4 = 5.4$$

$$q_{\pi}(s_3, \text{Study}) = 10 + 0 = 10$$

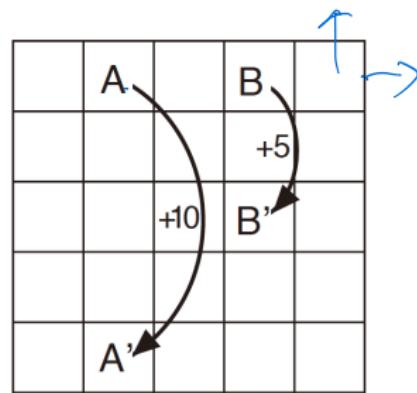
$$q_{\pi}(s_3, \text{Pub}) = 1 + 0.2V_1 + 0.4V_2 + 0.4V_3 = 1 + 0.2(-1.3) + 0.4(2.7) + 0.4(7.4) = 4.78$$

$$q_{\pi}(s_4, \text{Facebook}) = -1 + V_4 = -1 - 2.3 = -3.3$$

$$q_{\pi}(s_4, \text{Quit}) = 0 + V_1 = 0 - 1.3 = -1.3$$

Example: Bellman Expectation Equation in Student MDP

Example: Bellman Expectation Equation in Gridworld



(a)

try it!

(b)

What is the value function for the uniform random policy?

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

Bellman Expectation Equation (Matrix Form)

The Bellman expectation equation can be expressed concisely using the induced MRP,

$$v_\pi = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v_\pi$$

with direct solution

$$v_\pi = \underbrace{(I - \gamma \mathcal{P}^\pi)^{-1}}_{\text{underlined}} \mathcal{R}^\pi$$

Example: Bellman Expectation Equation (Matrix Form)

$$r=1; \quad \underline{R}^{\pi} = \sum_{a \in A} \pi(a|s) R_s^a, \quad \underline{R}^{\pi} = [R_s^{\pi}]$$

$$\underline{P}_{s,s'}^{\pi} = \sum_{a \in A} \pi(a|s) P_{s,s'}^a, \quad \underline{P}^{\pi} = [P_{s,s'}^{\pi}]$$

$$\underline{P}^{\pi} = \begin{bmatrix} S_1 & S_2 & S_3 & S_4 & \text{Sleep} \\ S_1 & 0.5 & & & \\ S_2 & & 0.5 & & \\ S_3 & 0.5 \times 0.2 & 0.5 \times 0.4 & 0.5 \times 0.4 & \\ S_4 & 0.5 & 0.5 & 0.5 & \\ \text{Sleep} & & & 0.5 & 1 \end{bmatrix}$$

$$\underline{V}^{\pi} = (I - r\underline{P}^{\pi})^{-1} \underline{R}^{\pi}$$

$$\underline{R}^{\pi} = \begin{bmatrix} S_1 & 0.5(-2) + 0.5(-1) \\ S_2 & 0.5(-2) + 0.5(0) \\ S_3 & 0.5(10) + 0.5(1) \\ S_4 & 0.5(-1) + 0.5(0) \\ \text{Sleep} & 0 \end{bmatrix} = \begin{bmatrix} -1.5 \\ -1 \\ 5.5 \\ -0.5 \\ 0 \end{bmatrix}$$

Example: Bellman Expectation Equation (Matrix Form)

Optimal Value Function

Definition

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

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

$$\pi^* = \arg \max_{\pi} v_{\pi}(s)$$

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

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) = \max_{\pi} [R_s^a + \gamma \sum_{s' \in S} P_{s,s'}^a \cdot v_{\pi}(s')]$$

$$= R_s^a + \gamma \sum_{s' \in S} P_{s,s'}^a \cdot \max_{\pi} v_{\pi}(s')$$

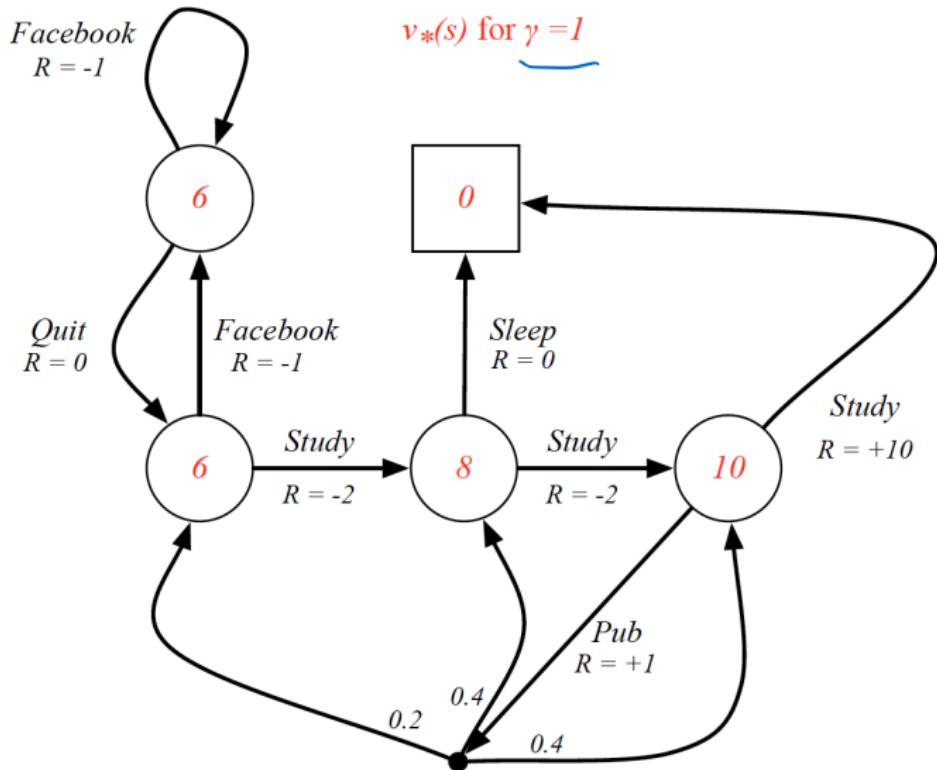
$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

$$(\pi^*)' = \arg \max_{\pi} q_{\pi}(s, a)$$

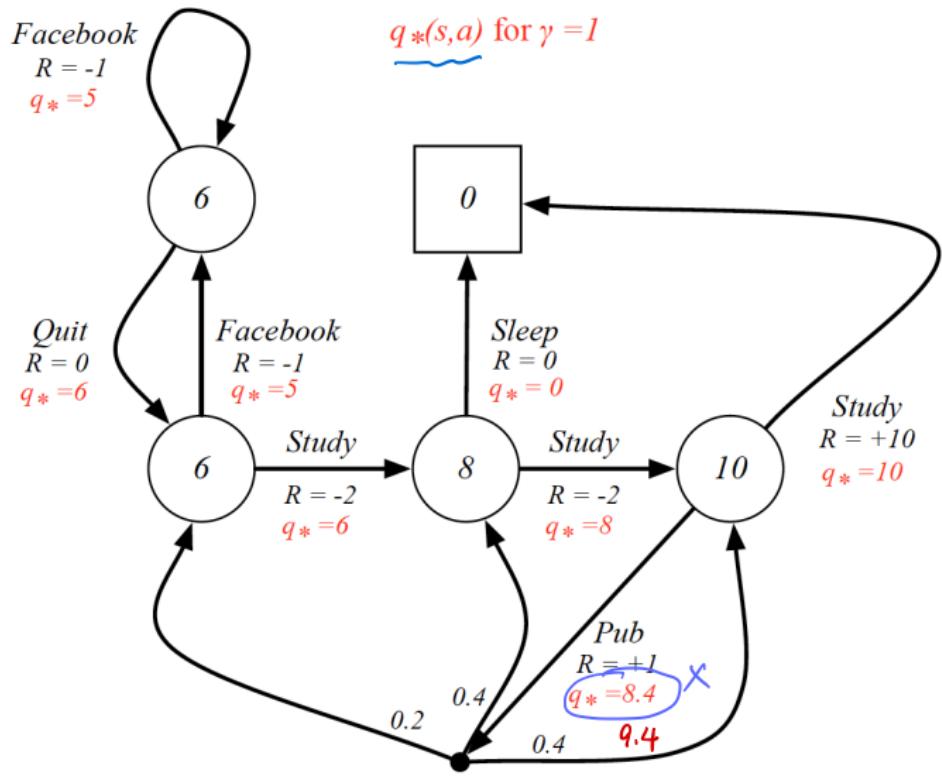
Intuitively $(\pi^*)' = \pi^*$.

- 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



Example: Optimal Action-Value Function for Student MDP

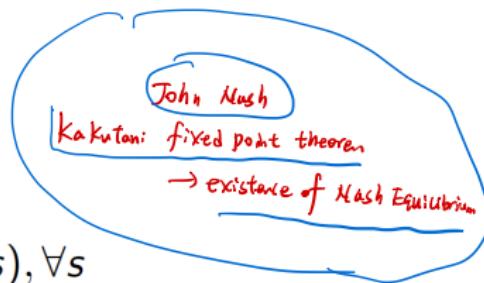


Optimal Policy

Define a partial ordering over policies

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

Banach fixed point theorem.



Theorem

For any Markov Decision Process

- There exists an optimal policy π_* 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

Linear programming (if $X_f > 0$ given constants)

$$\max_{f \in F} X_f$$

$$\text{s.t. } \sum_a p_f^a = 1$$

$$0 \leq p_f^a \leq 1$$

assume f^* is unique

$$\begin{aligned} \text{Solution } f^* &= \arg\max_{f \in F} X_f \\ p_f^* &= \begin{cases} 1 & \text{if } f = f^* \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

$$\begin{aligned} \text{if } \pi' \text{ is a random optimal policy.} \\ \text{or } \pi'(a|s) < 1, \forall a \in A. \\ V_{\pi}(s) &= V_{\pi'}(s) = \sum_{a \in A} \pi'(a|s) q_{\pi'}(s, a) \\ &\leq \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a) \\ &\leq \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a) \\ &= q_{\pi}(s, a^*) \quad \text{①} \end{aligned}$$

maximizing

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

$$\max_{f \in F} X_f = X_{\max}$$

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \arg\max_{a \in A} q_*(s, a) = a^* \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} V_{\pi}(s) &\geq V_{\pi'}(s) \\ &= \sum_{a \in A} \pi'(a|s) q_{\pi'}(s, a) \\ &= q_{\pi}(s, a^*) \quad \text{②} \end{aligned}$$

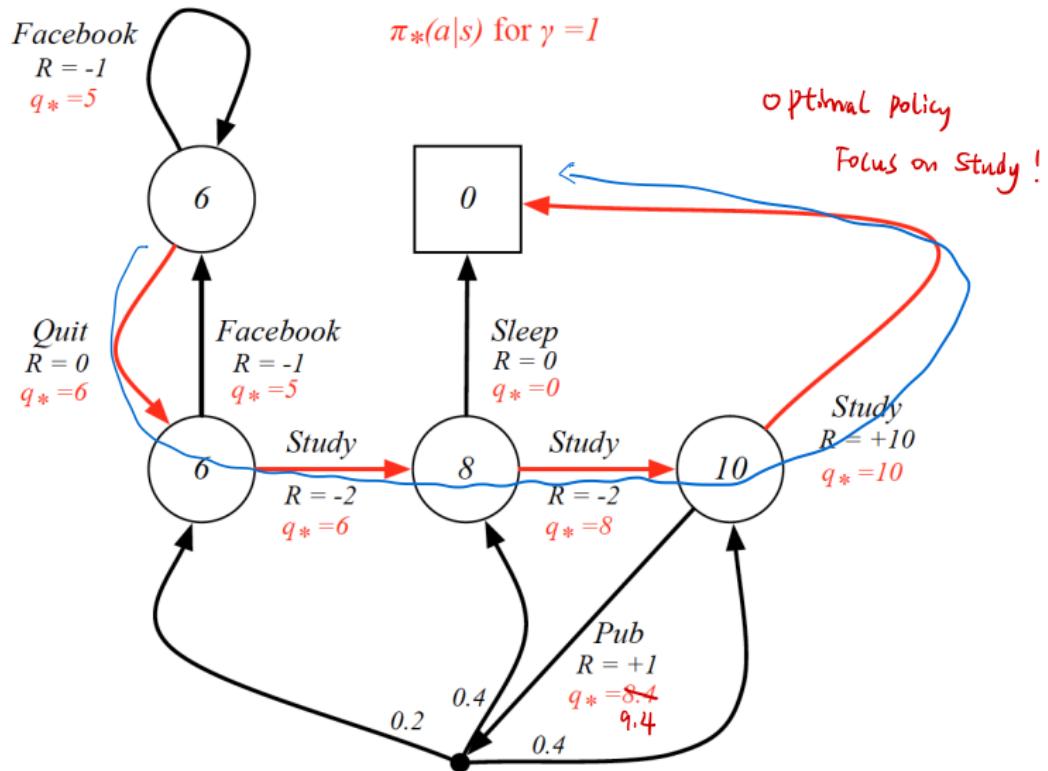
$\Rightarrow \pi(a|s) > 0, \forall a \in A$. random policy.

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

From ① we know if you have a choice between picking exactly the action you want (if $\pi(a|s) > 0$) versus picking a probability distribution over potentially optimal and non-optimal actions, you would always prefer to pick exactly the best action.

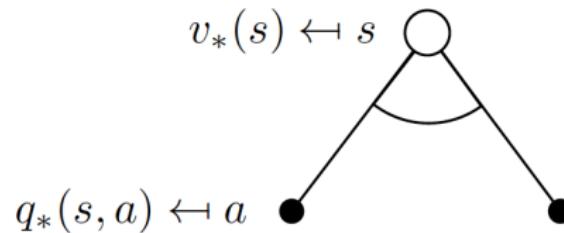
$$\text{From ① + ② } V_{\pi}(s) = q_{\pi}(s, a^*) = \max_{a \in A} q_{\pi}(s, a)$$

Example: Optimal Policy for Student MDP



Bellman Optimality Equation for v_*

The optimal value functions are recursively related by the Bellman optimality equations:



$$v_*(s) = \max_a q_*(s, a)$$

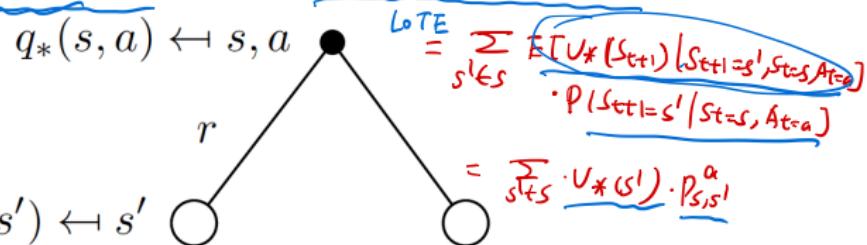
Bellman Optimality Equation for Q^*

$$1^{\circ} \quad q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{s,s'}^a v_{\pi}(s')$$

$$2^{\circ} \quad q^*(s, a) = \max_{\pi} q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{s,s'}^a \cdot \max_{\pi} v_{\pi}(s') = R_s^a + \gamma \sum_{s' \in S} p_{s,s'}^a \cdot v^*(s')$$

$$q^*(s, a) = \mathbb{E}[R_{t+1} + \gamma v^*(S_{t+1}) | S_t = s, A_t = a]$$

$$3^{\circ} \quad \mathbb{E}[R_{t+1} | S_t = s, A_t = a] = R_s^a, \quad \mathbb{E}[v^*(S_{t+1}) | S_t = s, A_t = a]$$



$$4^{\circ} \quad q^*(s, a) = \mathbb{E}[R_{t+1} + \gamma V^*(S_{t+1}) | S_t = s, A_t = a]$$

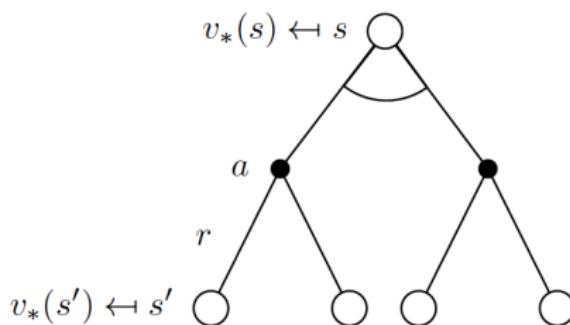
$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

$$q^*(s, a) = \mathbb{E}[R_{t+1} + \gamma V^*(S_{t+1}) | S_t = s, A_t = a]$$

Bellman Optimality Equation for V^*

$$V^*(s) = \max_a Q^*(s, a), \text{ and } Q^*(s, a) = E[R_{t+1} + \gamma V^*(S_{t+1}) | S_t=s, A_t=a]$$

$$v_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a]$$



$$Q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V^*(s')$$

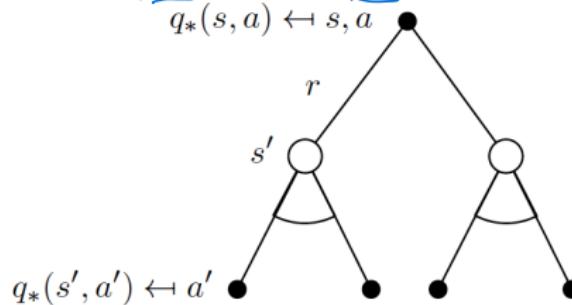
$$v_*(s) = \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s'))$$

Bellman Optimality Equation for Q^*

$$2. \mathbb{E} \left[\max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

LOTE $\sum_{s' \in S} \mathbb{E} \left[\max_{a'} q_*(S_{t+1}, a') \mid S_{t+1} = s', S_t = s, A_t = a \right] \cdot P(S_{t+1} = s' \mid S_t = s, A_t = a)$

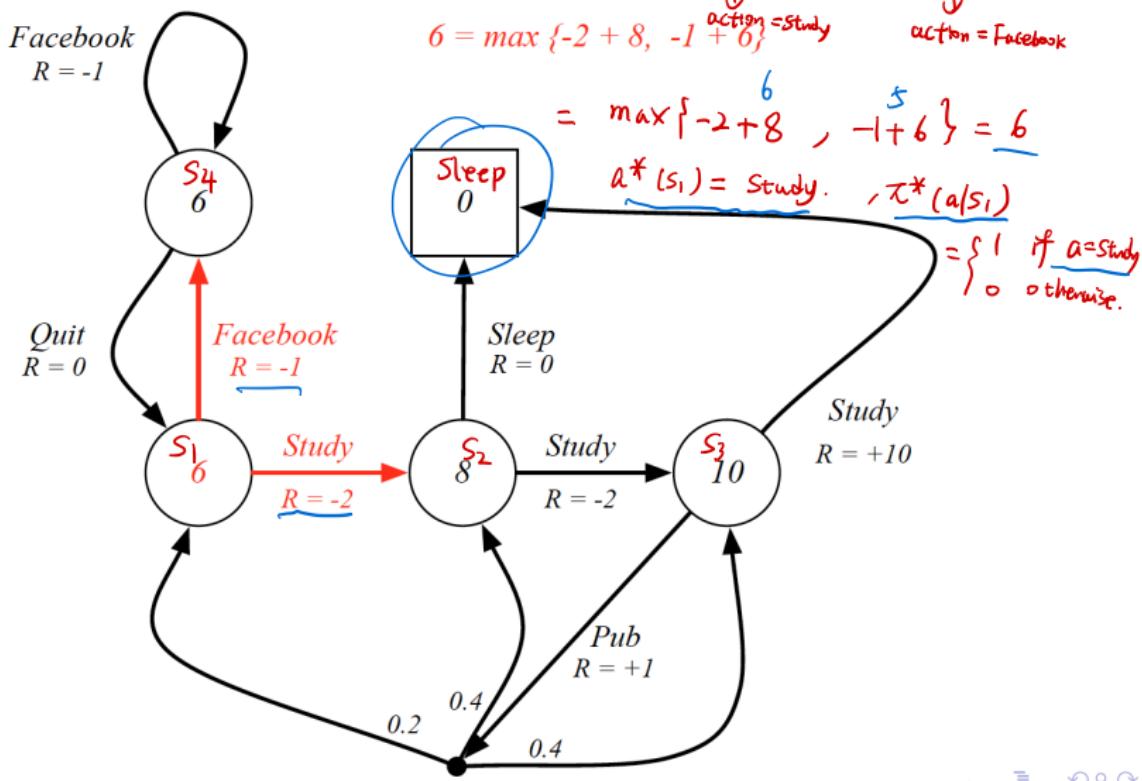
$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{s' \in S} \left(\max_{a'} q_*(s', a') \right) \cdot P_{s, s'}^a \end{aligned}$$



$$1. \underbrace{q_*(s, a)}_{\text{Bellman Optimality Equation}} = \mathcal{R}_s^a + \gamma \sum_{s' \in S} P_{s, s'}^a \underbrace{V_*(s')}_{\text{and } V_*(s') = \max_{a'} q_*(s', a')}$$

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in S} P_{s, s'}^a \max_{a'} q_*(s', a')$$

Example: Bellman Optimality Equation in Student MDP



Solution

Overall, we obtain $q_{\pi^*}(s, a)$ first, then obtain $V_{\pi^*}(s)$

$$1^{\text{st}}. \quad q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} p_{s, s'}^a V_{\pi}(s')$$

$$q_{\pi^*}(s, a) = \max_{\pi} q_{\pi}(s, a)$$

Since $V_{\pi}(\text{sleep}) \equiv 0, \forall \pi$, then we have

$$\begin{aligned} q_{\pi}(s_2, \text{sleep}) &= R_{s_2}^{\text{sleep}} + 1 \cdot V_{\pi}(\text{sleep}) = 0 + 0, \forall \pi \\ \Rightarrow q_{\pi^*}(s_2, \text{sleep}) &= 0 \quad \checkmark \end{aligned}$$

border condition.
terminal states

on the other hand,

$$\begin{aligned} q_{\pi}(s_3, \text{study}) &= R_{s_3}^{\text{study}} + 1 \cdot V_{\pi}(\text{sleep}) = 10 + 0 = 10, \forall \pi \\ \Rightarrow q_{\pi^*}(s_3, \text{study}) &= 10 \quad \checkmark \end{aligned}$$

Solution

$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{s,s'}^a \cdot \max_{a'} q^*(s', a')$$

$$2^0. \quad q^*(s_1, \text{study}) = R_{s_1}^{\text{Study}} + \gamma P_{s_1, s_2}^{\text{Study}} \max_{a'} q^*(s_2, a')$$

$$\begin{aligned} &= R_{s_1}^{\text{Study}} + \gamma P_{s_1, s_2}^{\text{Study}} \max(q^*(s_2, \text{sleep}), q^*(s_2, \text{study})) \\ &= -2 + \max[0, q^*(s_2, \text{study})] \end{aligned}$$

$$q^*(s_1, \text{facebook}) = R_{s_1}^{\text{Facebook}} + \gamma P_{s_1, s_4}^{\text{Facebook}} \cdot \max_{a'} q^*(s_4, a')$$

$$\begin{aligned} &= R_{s_1}^{\text{Facebook}} + \gamma P_{s_1, s_4}^{\text{Facebook}} \cdot \max[q^*(s_4, \text{Facebook}), q^*(s_4, \text{Quit})] \\ &= -1 + \max[q^*(s_4, \text{Facebook}), q^*(s_4, \text{Quit})] \end{aligned}$$

$$q^*(s_2, \text{sleep}) = 0$$

$$q^*(s_2, \text{study}) = -2 + \max[0, q^*(s_3, \text{Pub})]$$

$$q^*(s_3, \text{study}) = 10$$

Solution

$$\begin{aligned} \underline{q^*(s_3, \text{Pub})} &= 1 + 0.2 \cdot \max [q^*(s_1, \text{Study}), q^*(s_1, \text{Facebook})] \\ &\quad + 0.4 \max [q^*(s_2, \text{Study}), 0] \\ &\quad + 0.4 \max [q^*(s_3, \text{Pub}), 10] \end{aligned}$$

$$\underline{q^*(s_4, \text{Facebook})} = -1 + \max [q^*(s_4, \text{Facebook}), q^*(s_4, \text{Quit})]$$

$$\underline{q^*(s_4, \text{Quit})} = 0 + \max [q^*(s_1, \text{Study}), q^*(s_1, \text{Facebook})]$$

Solution

$$3^0. \quad \underbrace{q*(S_2, \text{study})}_{= -2 + \max[10, q*(S_3, \text{pub})]} \geq 8$$

$$\begin{aligned} \underbrace{q*(S_1, \text{study})}_{=} & -2 + \max[0, q*(S_2, \text{study})] \geq 6 \\ & = -2 + \underbrace{q*(S_2, \text{study})}_{=} \end{aligned}$$

$$\begin{aligned} \underbrace{q*(S_4, \text{Facebook})}_{?} & = -1 + \max[\underbrace{q*(S_4, \text{Facebook})}_{}, q*(S_4, \text{Quit})] \\ & = -1 + \underbrace{q*(S_4, \text{Facebook})}_{\text{X}} \\ & = -1 + \underbrace{q*(S_4, \text{Quit})}_{\text{X}} < \underbrace{q*(S_4, \text{Quit})}_{\text{X}} \end{aligned}$$

$$\begin{aligned} \underbrace{q*(S_1, \text{Facebook})}_{=} & -1 + \max[q*(S_4, \text{Facebook}), q*(S_4, \text{Quit})] \\ & = -1 + \underbrace{q*(S_4, \text{Quit})}_{=} \end{aligned}$$

Solution

$$q*(s_4, \text{quit}) = \max [q*(s_1, \text{study}), q*(s_1, \text{Facebook})]$$

$$= \max [\underbrace{q*(s_1, \text{study})}_{?}, \underbrace{q*(s_4, \text{quit}) - 1}_{\times}]$$

$$= \underbrace{q*(s_4, \text{quit}) - 1}_{\times}$$

$$= q*(s_1, \text{study}) = -2 + \underbrace{q*(s_2, \text{study})}_{}$$

$$q*(s_3, \text{Pub}) = 1 + 0.2 \max [\underbrace{q*(s_4, \text{quit})}_{q*(s_1, \text{study})}, \underbrace{-1 + q*(s_4, \text{quit})}_{q*(s_1, \text{Facebook})}]$$

$$+ 0.4 \max [\underbrace{q*(s_2, \text{study})}_{\geq 8}, 0]$$

$$+ 0.4 \max [q*(s_3, \text{Pub}), 10]$$

$$= 1 + 0.2 q*(s_4, \text{quit}) + 0.4 q*(s_2, \text{study}) + 0.4 \max [q*(s_3, \text{Pub}), 10]$$

$$= 1 + 0.2 [-2 + q*(s_2, \text{study})] + 0.4 q*(s_2, \text{study}) + 0.4 \max [q*(s_3, \text{Pub}), 10]$$

$$= 0.6 + 0.6 q*(s_2, \text{study}) + 0.4 \max [q*(s_3, \text{Pub}), 10]$$

Solution

4°. if $\max[q_*(s_3, \text{pub}), 10] = q_*(s_3, \text{pub}) (\geq 10)$

$$\Rightarrow q_*(s_3, \text{Pub}) = 0.6 + 0.6 q_*(s_2, \text{Study}) + 0.4 q_*(s_3, \text{Pub})$$

$$\Rightarrow q_*(s_3, \text{Pub}) = 1 + q_*(s_2, \text{Study})$$

On the other hand, we have

$$q_*(s_2, \text{Study}) = -2 + \max[10, q_*(s_3, \text{Pub})] = -2 + q_*(s_3, \text{Pub})$$

$$\Rightarrow q_*(s_3, \text{Pub}) = -1 + q_*(s_3, \text{Pub})$$

Thus $\underline{q_*(s_3, \text{Pub}) < 10}$, $\underline{\max[q_*(s_3, \text{Pub}), 10]} = 10$

Solution

5°. Then we have

$$\underline{q^*(s_3, \text{pub})} = 0.6 + 0.6 \underline{q^*(s_2, \text{study})} + 4 = 4.6 + 0.6 \underline{q^*(s_2, \text{study})}$$

On the other hand, $\underline{q^*(s_2, \text{study})} = -2 + \max [10, \underline{q^*(s_3, \text{pub})}]$
 $= -2 + 10 = 8$

$\Rightarrow \underline{q^*(s_1, \text{Study})} = -2 + \underline{q^*(s_2, \text{study})} = -2 + 8 = 6$

$$\underline{q^*(s_3, \text{Pub})} = 4.6 + 0.6 \underline{q^*(s_2, \text{study})} = 4.6 + 0.6 \times 8 = 9.4$$

$$\underline{q^*(s_4, \text{Quit})} = \underline{q^*(s_1, \text{Study})} = 6$$

$$\underline{q^*(s_4, \text{Facebook})} = -1 + \underline{q^*(s_4, \text{Quit})} = -1 + 6 = 5$$

$$\underline{q^*(s_1, \text{Facebook})} = -1 + \underline{q^*(s_4, \text{Quit})} = -1 + 6 = 5$$

Solution

$$6^\circ. \quad \underline{V^*(s) = \max_a q^*(s, a)}$$

$$\Rightarrow \underline{V^*(s_1) = \max_a q^*(s_1, a)} = \max [q^*(s_1, \text{Study}), q^*(s_1, \text{Facebook})] \\ = \max [6, 5] = 6$$

$$\underline{V^*(s_2) = \max [q^*(s_2, \text{Study}), q^*(s_2, \text{Sleep})]} \\ = \max [8, 0] = 8$$

$$\underline{V^*(s_3) = \max [q^*(s_3, \text{Study}), q^*(s_3, \text{Pub})]} \\ = \max [10, 9.4] = 10$$

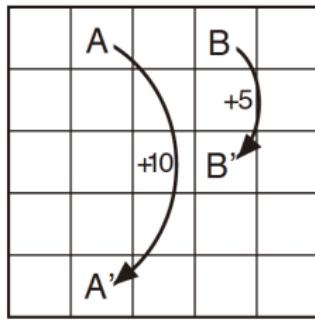
$$\underline{V^*(s_4) = \max [q^*(s_4, \text{Dmit}), q^*(s_4, \text{Facebook})]} \\ = \max [6, 5] = 6$$

$$\underline{V^*(\text{Sleep}) = 0}$$

Solution

Solution

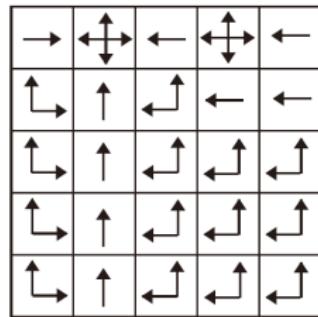
Example: Bellman Optimality Equation in Gridworld



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b) v_*



c) π_*

What is the optimal value function over all possible policies?
What is the optimal policy?

try it !

Solution

Solving the Bellman Optimality Equation in General

- Bellman Optimality Equation is non-linear
 - No closed form solution (in general)
 - Many iterative solution methods
 - ▶ Value Iteration
 - ▶ Policy Iteration
 - ▶ Q-learning
 - ▶ Sarsa
- Control. LQR*

Solving the Bellman Optimality Equation

- Finding an optimal policy by solving the Bellman Optimality Equation requires the following:
 - ▶ accurate knowledge of environment dynamics
 - ▶ we have enough space and time to do the computation
 - ▶ the Markov Property
- How much space and time do we need?
 - ▶ polynomial in number of states
 - ▶ BUT, number of states is often huge
 - ▶ So exhaustive sweeps of the state space are not possible

$$P_{s,s'}^a$$

Approximation and Reinforcement Learning

- RL methods: Approximating Bellman optimality equations
 - Balancing reward accumulation and system identification (model learning) in case of unknown dynamics *Unknown: $p_{s,s'}$*
 - The on-line nature of reinforcement learning makes it possible to approximate optimal policies in ways that put more effort into learning to make good decisions for frequently encountered states, at the expense of less effort for infrequently encountered states.
- Exploration - Exploitation Tradeoff

Outline

1 Introduction

2 Markov Reward Process

3 Markov Decision Process

4 References

Main References

- Reinforcement Learning: An Introduction (second edition), R. Sutton & A. Barto, 2018.
- RL course slides from Richard Sutton, University of Alberta.
- RL course slides from David Silver, University College London.