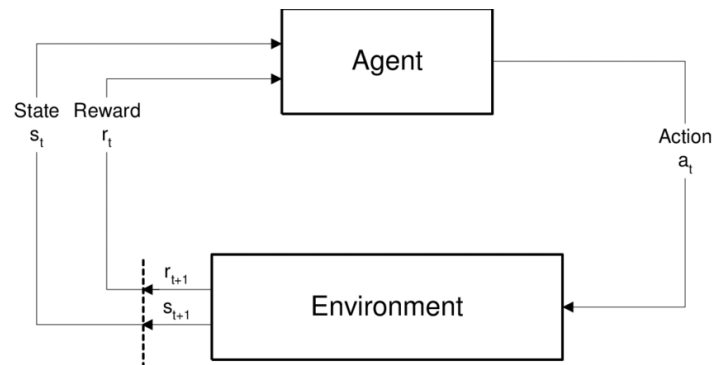


Last time ...



$S_0, A_0, R_1, S_1, A_1, R_2, \dots$   
trajectory

- Discounted Future rewards

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

$\gamma$  discount rate  $\in [0, 1]$

- value function

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

Intuitively: given a state  $s$ , tells me the expected discounted future rewards if I follow policy  $\pi$ .

- action-value function

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

Intuitively: given a state  $s$  AND action  $a$ , tells me the expected discounted future rewards if I follow policy  $\pi$ .

Recall MDP (Markov Decision Process).

Formally: MDP is a tuple  $\langle S, A, P, R, \gamma \rangle$

$S$ : finite set of states

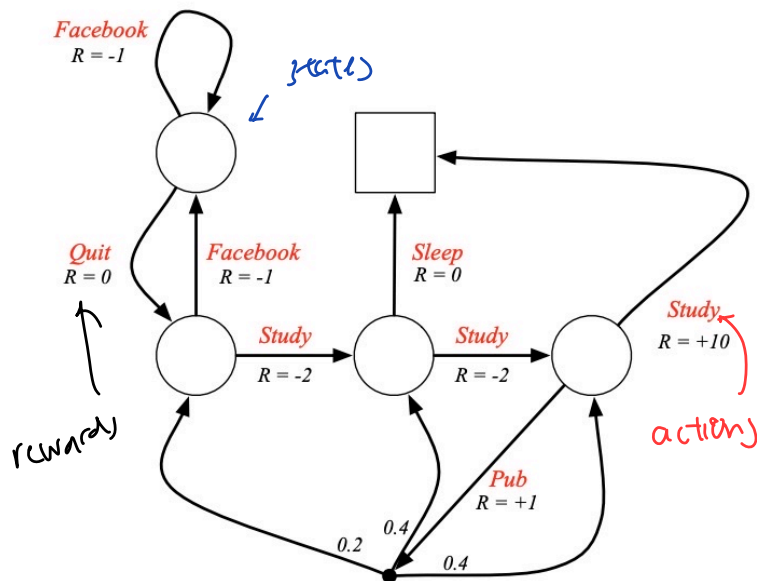
$A$ : finite set of actions

$P$ : transition probability b/w. states:  $P(S_{t+1} = s' | S_t = s, A_t = a_t)$

$R$ : Reward function:  $\mathbb{R} [R_{t+1} | S_t = s, A_t = a_t]$

$\gamma$ : discount factor  $\gamma \in [0, 1]$

Example:



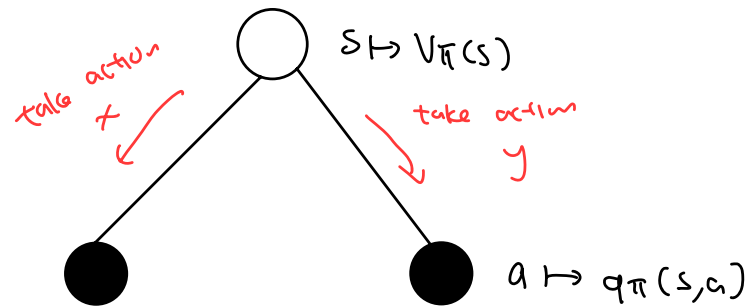
Student MDP

State-value function:

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

Can be expressed as

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

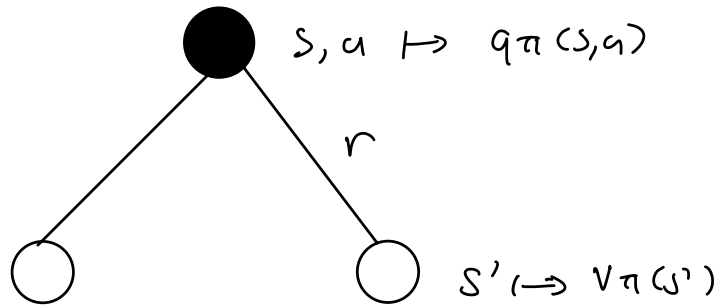


$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

taking the average across all actions.

Similarly for  $q_\pi(s, a) \dots$

$$\begin{aligned}
 q_\pi(s, a) &= \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a] + \sum_{s' \in \mathcal{S}} P(s' \mid s, a) V_\pi(s') \\
 &= \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} P_{ss'}^a V_\pi(s')
 \end{aligned}$$



Recall, the Bellman Equation for  $v_\pi$  is given by

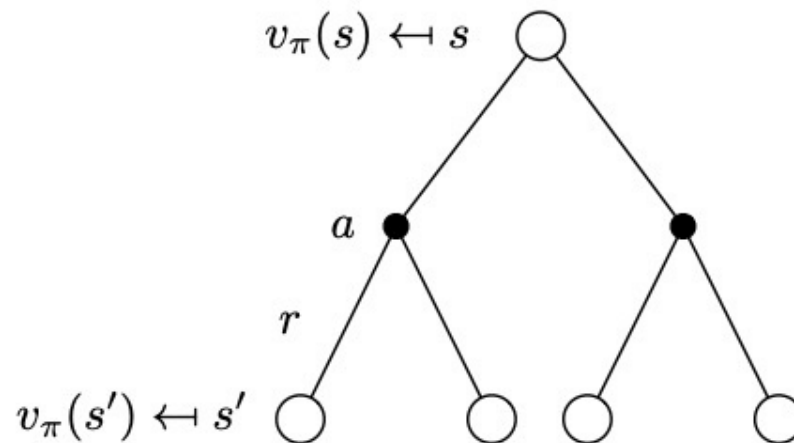
$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma v_\pi(s')] \quad \forall s \in S$$

- we sum over  $a, s', r$
- can be treated as an expected value
- for each triple  $(a, s', r)$ , we compute its probability  $\pi(a|s) p(s', r|s, a)$ 
  - weighted by  $[r + \gamma v_\pi(s')]$
  - Sum over all possibilities → expected value

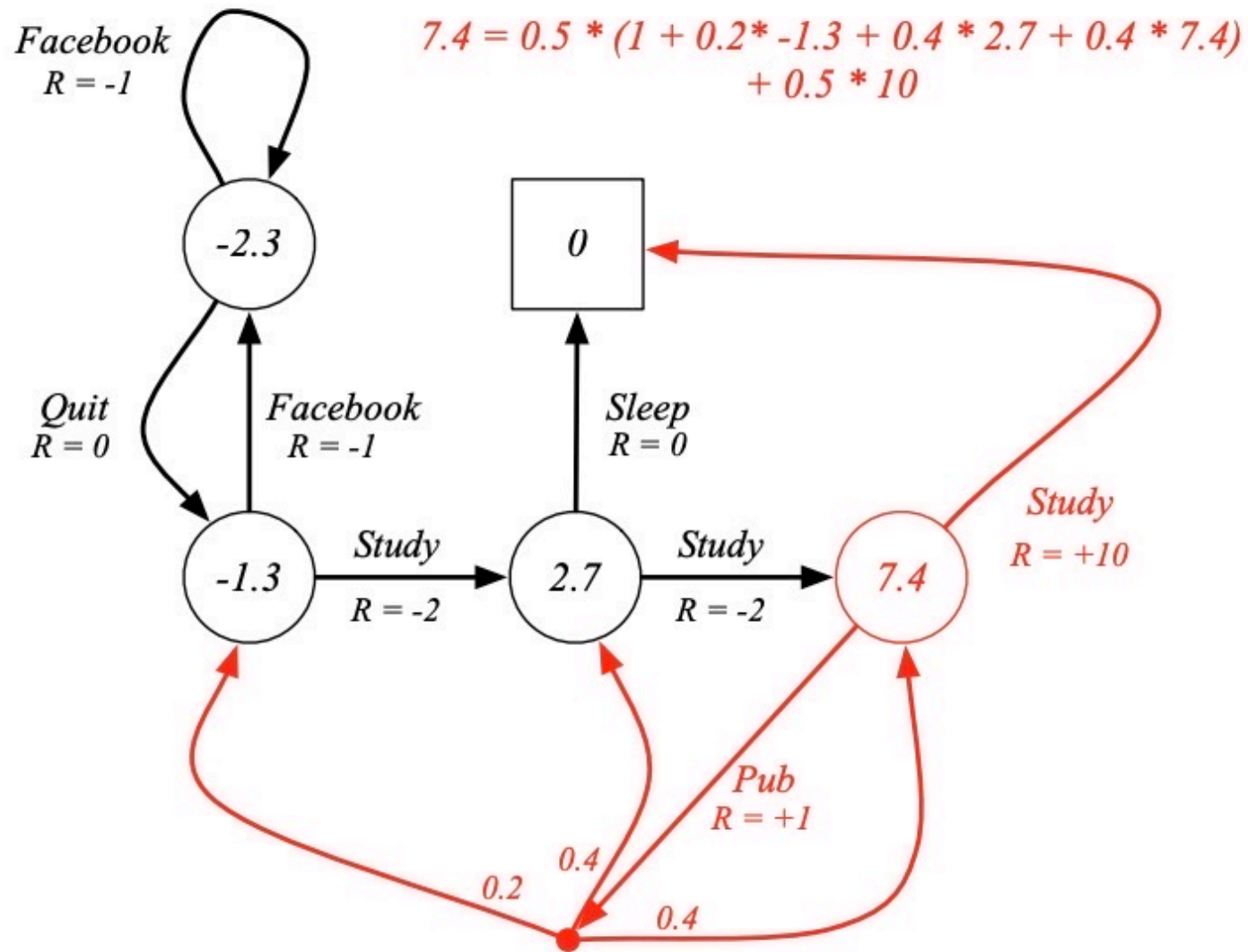
② expresses the relationship between the value of a state and value of successor states.

Intuitively: Value of state = discounted value of expected next state + reward along the way

To solve a reinforcement problem = finding  $V_\pi^*(s) = \max_\pi V_\pi(s)$  and  $q_\pi^*(s, a) = \max_\pi q_\pi(s, a)$



# Example



To solve for the  $V^*(s)$  and  $q^*(s,a)$ , we can use:

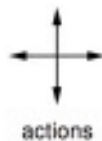
- value iteration
- policy iteration
- Q-learning
- Sarsa

Today: DP & policy evaluation / iteration

### Dynamic Programming

DP: Algorithms used to find optimal policies which have complete knowledge  $(p(s',r|s,a))$ . → not realistic assumption

Example:



actions

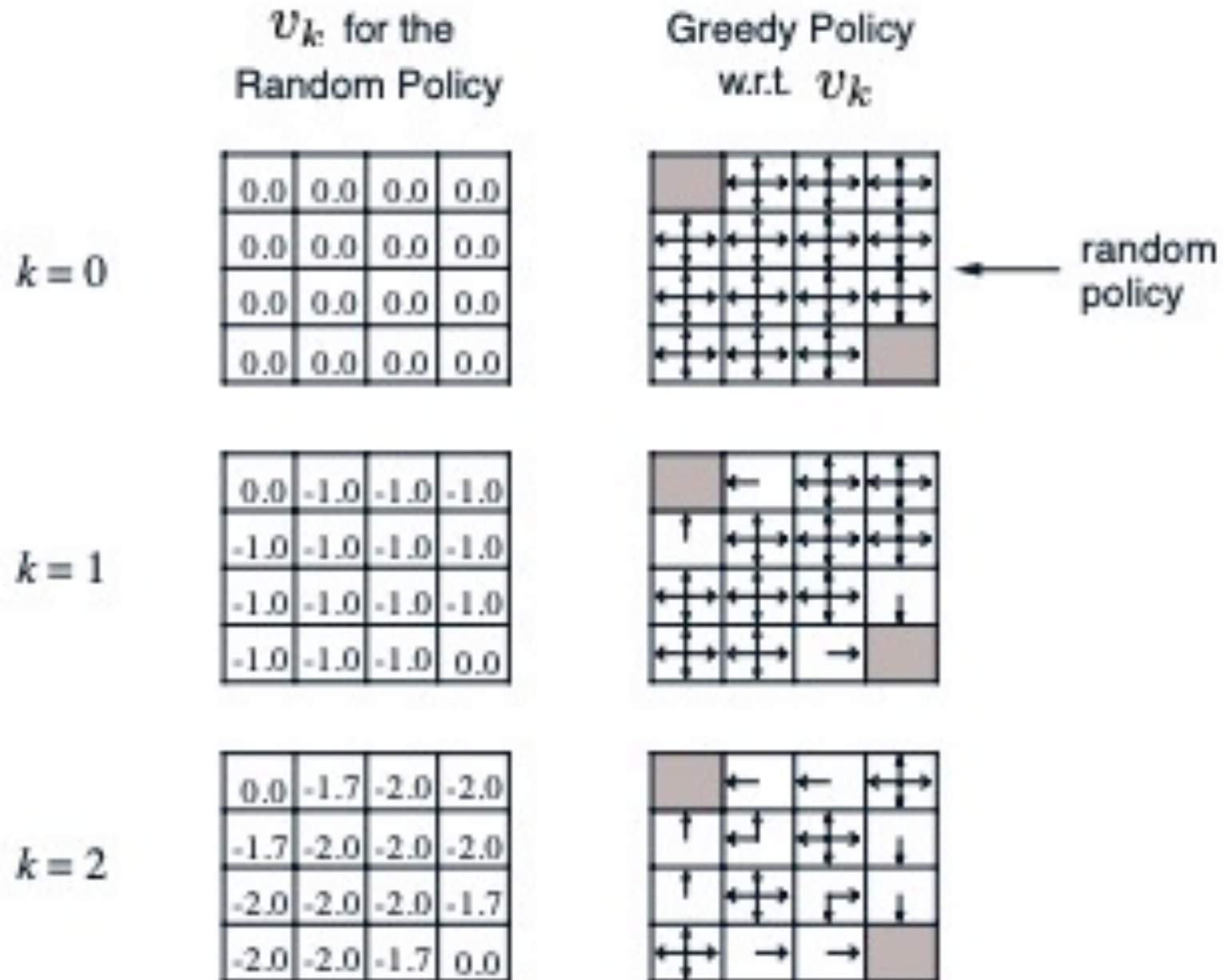
	1	2	
4	5	6	7
8	9	10	11
12	13	14	

$r = -1$   
on all transitions

- Undiscounted episodic MDP ( $\gamma = 1$ )
- Nonterminal states 1, ..., 14
- One terminal state (shown twice as shaded squares)
- Actions leading out of the grid leave state unchanged
- Reward is  $-1$  until the terminal state is reached
- Agent follows uniform random policy

$$\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$$

# POLICY ITERATION





$k = 3$ 

0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

 $k = 10$ 

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

 $k = \infty$ 

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



optimal  
policy

How do we improve a policy?

A: Given a policy  $\pi$ ,

① Evaluate the policy

$$V_{\pi}(s) = \mathbb{E}[G_t | s_t = s]$$

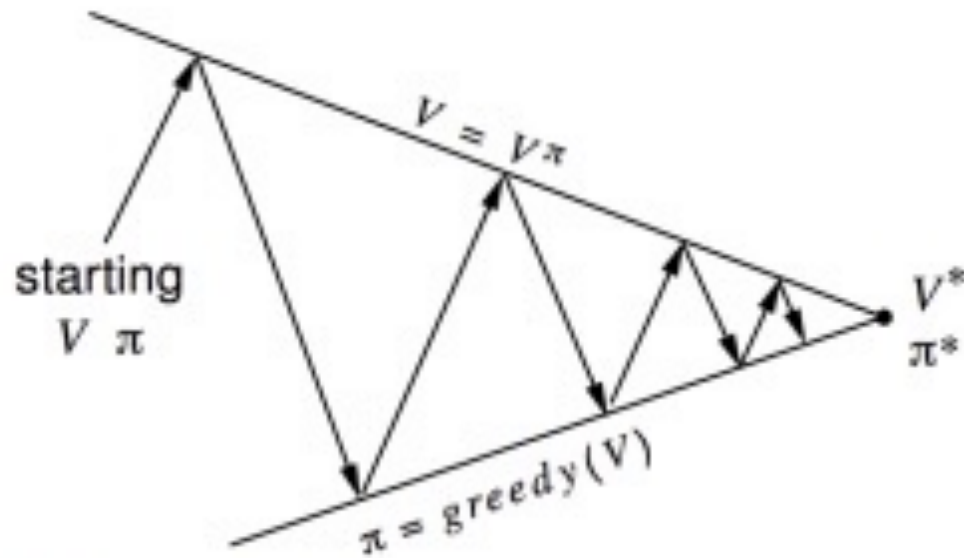
② Improve policy by acting greedily wrt.  $V_{\pi}$

$$\pi' = \text{greedy}(V_{\pi})$$

What is greedy?  $\rightarrow$  taking the maximum rewards / value

NOTE: POLICY ITERATION ALWAYS converges to  $\pi^*$

In general...

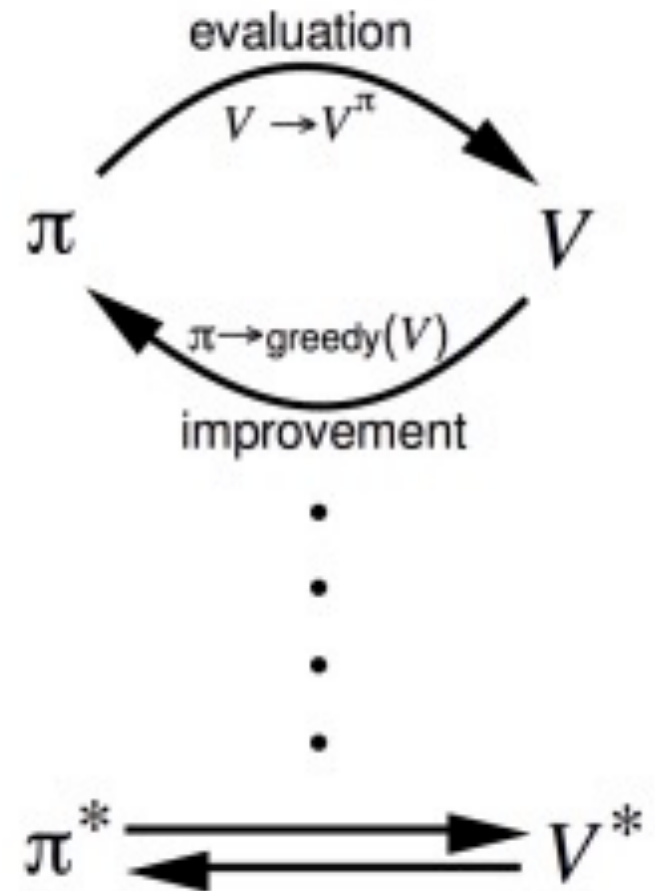


**Policy evaluation** Estimate  $v_\pi$

**Any** policy evaluation algorithm

**Policy improvement** Generate  $\pi' \geq \pi$

**Any** policy improvement algorithm



Next time:

- Value iteration
- Code in notebook
- Start thinking about projects