

Dynamic Programming

- Sutton, Richard S., and Barto, Andrew G. Reinforcement learning: An introduction. MIT press, 2018.
 - <http://www.incompleteideas.net/book/the-book-2nd.html>
 - Chapter 4
- Puterman, Martin L. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
 - Chapter 4
- David Silver lecture on Dynamic Programming
 - <https://www.youtube.com/watch?v=Nd1-UUMVfz4>

- A classical algorithm for computing solutions to problems that can be separated into subproblems with known solutions
 - E.g., all shortest paths
 - Essentially, store all subproblem solutions in a table and reuse them when necessary
- If we have a finite (state and action) MDP, we can find the optimal policy by incremental search
 - Find the optimal policy for 1 step, then 2 steps, etc.
 - Actually done backwards in time
- Polynomial complexity in the number of states
 - Number of states can be large

- In order to find the best policy, we first need a way to evaluate policies
 - i.e., compute the state-value function $v_{\pi}(s)$ for each state s
 - Also compute the action-value function $q_{\pi}(s, a)$
- Every time we change the policy, we need to evaluate it (in order to check if we improved it)
- So far, we've seen one way to compute state values
 - How?
 - For a given policy π , compute the matrix form of the value function vector

- In the infinite-horizon case, we can use the Bellman equation:

$$v_{\pi}(s) = R_{\pi}(s) + \gamma \sum_{a,s'} P(s, a, s') \pi(a|s) v_{\pi}(s')$$

- which can be rewritten in matrix form:

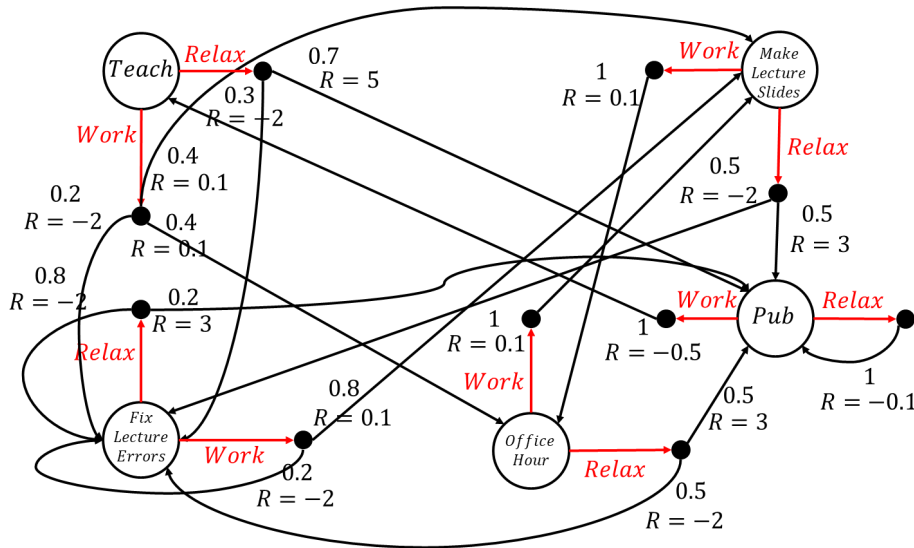
$$v_{\pi}(\mathbf{s}) = R_{\pi}(\mathbf{s}) + \gamma \mathbf{P}_{\pi} v_{\pi}(\mathbf{s})$$

- Thus, the state value vector is:

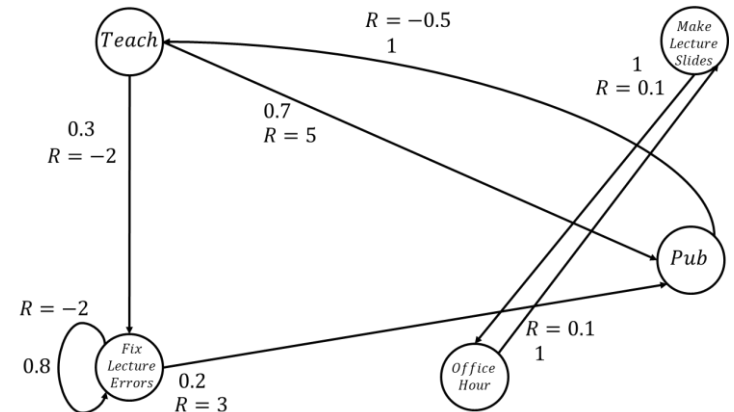
$$v_{\pi}(\mathbf{s}) = (\mathbf{I} - \gamma \mathbf{P}_{\pi})^{-1} R_{\pi}(\mathbf{s})$$

Workday example, MDP -> MRP

MDP



MRP

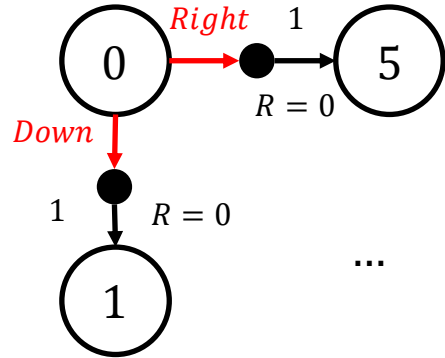
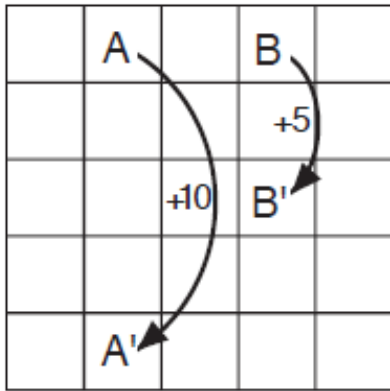


Policy

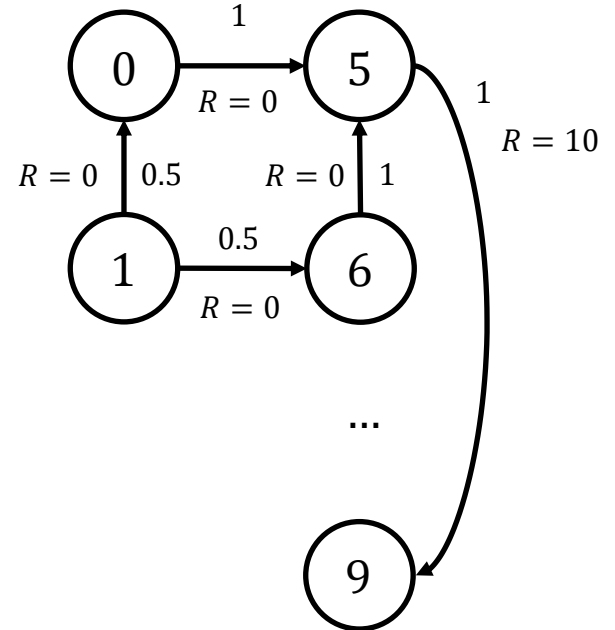
- Let's define π as follows:
 - $\pi(\text{Teach}) = \text{Relax}$
 - $\pi(\text{OH}) = \text{Work}$
 - $\pi(\text{MLS}) = \text{Work}$
 - $\pi(\text{FLE}) = \text{Relax}$
 - $\pi(\text{Pub}) = \text{Work}$

Gridworld example, MDP \rightarrow MRP

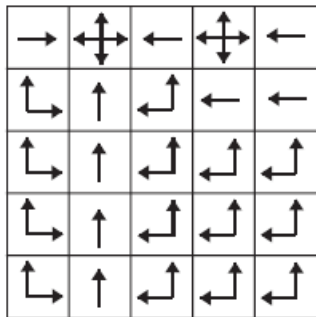
MDP



MRP



Policy



- Inverting $I - \gamma P$ may be expensive if number of states is large
- Another approach is to use linear systems theory!
- Start from a random initialization for $v(\mathbf{s}) = v_0(\mathbf{s})$
- Look at linear system

$$v_k(\mathbf{s}) = R(\mathbf{s}) + \gamma P v_{k-1}(\mathbf{s})$$

- System is stable. Why?
 - All entries (and eigenvalues) of γP are < 1 (for $\gamma < 1$)
- System converges to unique solution $(I - \gamma P)^{-1}R(\mathbf{s})$
 - Why?
 - If $v_{k+1}(\mathbf{s}) = v_k(\mathbf{s}) = \mathbf{v}$, then $\mathbf{v} = R(\mathbf{s}) + \gamma P \mathbf{v}$, i.e.,
$$\mathbf{v} = (I - \gamma P)^{-1}R(\mathbf{s})$$
 - Keep in mind $I - \gamma P$ is not invertible when $\gamma = 1$

- Recall state values are (for $\gamma = 0.9$)

$$(I - \gamma P)^{-1}R(\mathbf{s}) = [5.54 \quad 1 \quad 1 \quad -0.69 \quad 4.49]^T$$

- Using iterative evaluation

– Starting with $\mathbf{v}_0 = [0 \quad 0 \quad 0 \quad 0 \quad 0]^T$

$$\mathbf{v}_{10} = [4.59 \quad 0.65 \quad 0.65 \quad -1.58 \quad 3.64]^T$$

$$\mathbf{v}_{30} = [5.43 \quad 0.96 \quad 0.96 \quad 0 - 0.80 \quad 4.38]^T$$

$$\mathbf{v}_{50} = [5.53 \quad 0.99 \quad 0.99 \quad 0 - 0.70 \quad 4.47]^T$$

- After 50 iterations, converged within 0.01 if the true values

Gridworld Example, iterative value evolution

- Recall state values are
- Using iterative evaluation
 - Starting from $v_0 = \mathbf{0} \in \mathbb{R}^{25}$

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

- v_{10} is

14.31	15.90	14.31	10.90	9.81
12.88	14.31	12.88	11.59	10.44
11.59	12.88	11.59	10.44	5.90
10.44	11.59	10.44	5.90	5.31
5.90	10.44	5.90	5.31	4.78

- v_{50} is

21.86	24.29	21.86	19.29	17.36
19.68	21.86	19.68	17.71	15.94
17.71	19.68	17.71	15.94	14.29
15.94	17.71	15.94	14.29	12.86
14.29	15.94	14.29	12.86	11.58

- After 50 iterations, converge within 0.15 of true values
- In general, can stop iterating when $\|v_{k+1} - v_k\| \leq \epsilon$
 - Where ϵ is a hyperparameter

What about the finite-horizon case?

- Evaluate a policy recursively, starting from the last step T
 - Use the Bellman equation

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

- Remember, the policy and the value function may be time-dependent in the finite-horizon case!

- A policy π is better than another policy π' if

$$v_{\pi}(s) \geq v_{\pi'}(s), \forall s \in S$$

- A policy π^* is optimal if there exists no better policy than π^*
- Turns out the optimal policy also has a nice recursive property

$$\pi_*(s) = \mathit{arg} \max_a q_{\pi_*}(s, a)$$

- Another version of the Bellman equation
- Pick the action with the highest value
- Obvious in a sense
 - But any policy that satisfies the Bellman optimality equation is optimal

- Bellman optimality equation:

$$\begin{aligned}\pi_*(s) &= \mathit{arg} \max_a q_{\pi_*}(s, a) \\ &= \mathit{arg} \max_a \mathbb{E}_* [R_{t+1} + \gamma v_{\pi_*}(S_{t+1}) | S_t = s, A_t = a]\end{aligned}$$

- Proof by induction backward in time (for finite horizon T):

– Base case ($t = T - 1$):

- Only one step to make
- For any state s , the optimal action is

$$\begin{aligned}\mathit{arg} \max_a \mathbb{E}_* [R_T | S_{T-1} = s, A_{T-1} = a] &= \\ &= \mathit{arg} \max_a q_{\pi_*}(s, a)\end{aligned}$$

- So π_* is optimal by construction

- Proof by induction backward in time (for finite horizon T):
 - Inductive case:
 - Assume π_* is optimal at time $t + 1$
 - i.e., $v_{\pi_*}^{t+1}(s) \geq v_{\pi'}^{t+1}(s), \forall s, \pi'$
 - What do we need to show?

$$v_{\pi_*}^t(s) \geq v_{\pi'}^t(s), \forall s, \pi'$$

- Using the Bellman equation for π_* :

$$\begin{aligned} v_{\pi_*}^t(s) &= q_{\pi_*}^t(s, \pi_*(s)) \\ &= \max_a [q_{\pi_*}^t(s, a)] \\ &= \max_a \left[\mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}^{t+1}(S_{t+1}) | S_t = s, A_t = a] \right] \end{aligned}$$

- Proof by induction backward in time (for finite horizon T):

- Assume π_* is optimal at time $t + 1$

- i.e., $v_{\pi_*}^{t+1}(s) \geq v_{\pi'}^{t+1}(s), \forall s, \pi'$

- Consider any other policy π'

$$\begin{aligned} v_{\pi_*}^t(s) &= \max_a \left[\mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}^{t+1}(S_{t+1}) | S_t = s, A_t = a] \right] \\ &= \max_a \left[R_e(s, a) + \sum_{s'} \gamma v_{\pi_*}^{t+1}(s') P(s, a, s') \right] \\ &\geq \max_a \left[R_e(s, a) + \sum_{s'} \gamma v_{\pi'}^{t+1}(s') P(s, a, s') \right] \end{aligned}$$

- Inequality true for any a , so true for max also

- Proof by induction backward in time (for finite horizon T):

$$\begin{aligned} v_{\pi_*}^t(s) &= \max_a \left[\mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}^{t+1}(S_{t+1}) | S_t = s, A_t = a] \right] \\ &= \max_a \left[R(s, a) + \sum_{s'} \gamma v_{\pi_*}^{t+1}(s') P(s, a, s') \right] \\ &\geq \max_a \left[R(s, a) + \sum_{s'} \gamma v_{\pi'}^{t+1}(s') P(s, a, s') \right] \\ &\geq R(s, \pi'(s)) + \sum_{s'} \gamma v_{\pi'}^{t+1}(s') P(s, \pi'(s), s') \\ &= \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_{\pi'}^{t+1}(S_{t+1}) | S_t = s, A_t = \pi'(s)] \\ &= \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_{\pi'}^{t+1}(S_{t+1}) | S_t = s] \\ &= v_{\pi'}^t(s) \end{aligned}$$

Deterministic
policy version

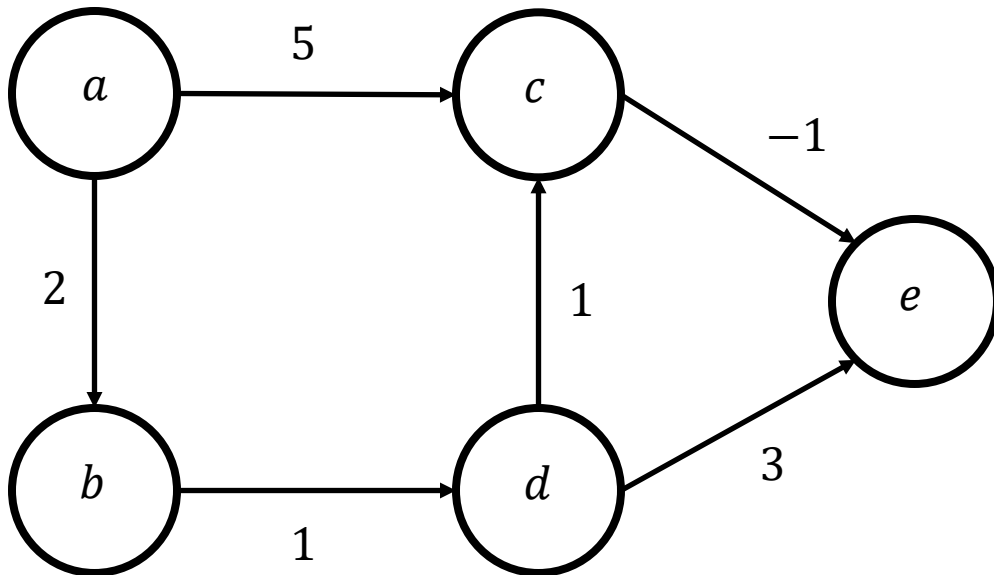
- Proof by induction backward in time (for finite horizon T):

$$\begin{aligned} v_{\pi_*}^t(s) &= \max_a \left[\mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}^{t+1}(S_{t+1}) | S_t = s, A_t = a] \right] \\ &= \max_a \left[R(s, a) + \sum_{s'} \gamma v_{\pi_*}^{t+1}(s') P(s, a, s') \right] \\ &\geq \max_a \left[R(s, a) + \sum_{s'} \gamma v_{\pi'}^{t+1}(s') P(s, a, s') \right] \\ &\geq \sum_{a'} \pi'(a'|s) \left[R(s, a') + \sum_{s'} \gamma v_{\pi'}^{t+1}(s') P(s, a', s') \right] \\ &= v_{\pi'}^t(s) \end{aligned}$$

Stochastic
policy
version

Marginalize
over all
actions

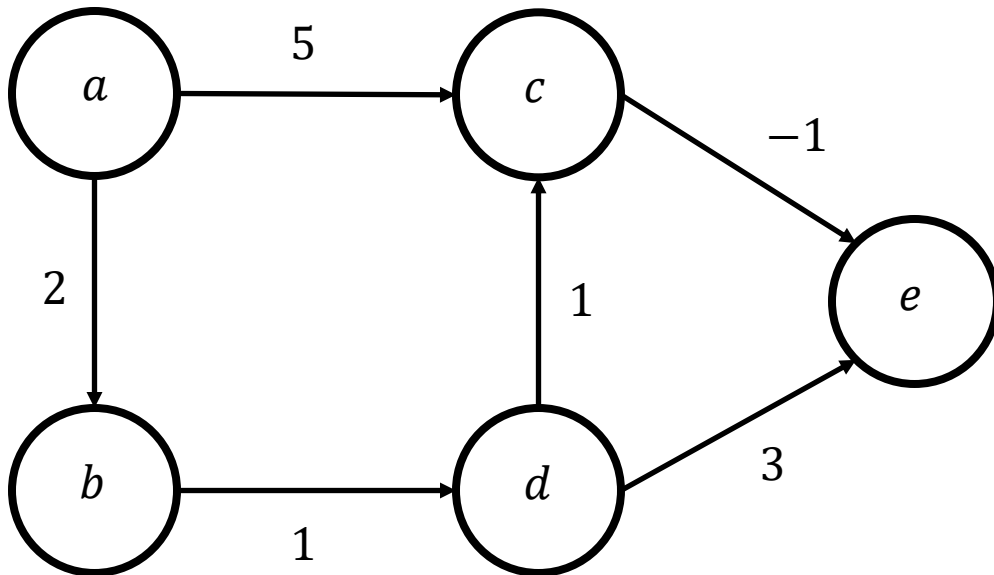
- A famous application of dynamic programming
 - Compute shortest paths from a node to all other nodes in a graph
 - E.g., all shortest paths from a to other nodes



1 step

a	0		
b	2		
c	5		
d	∞		
e	∞		

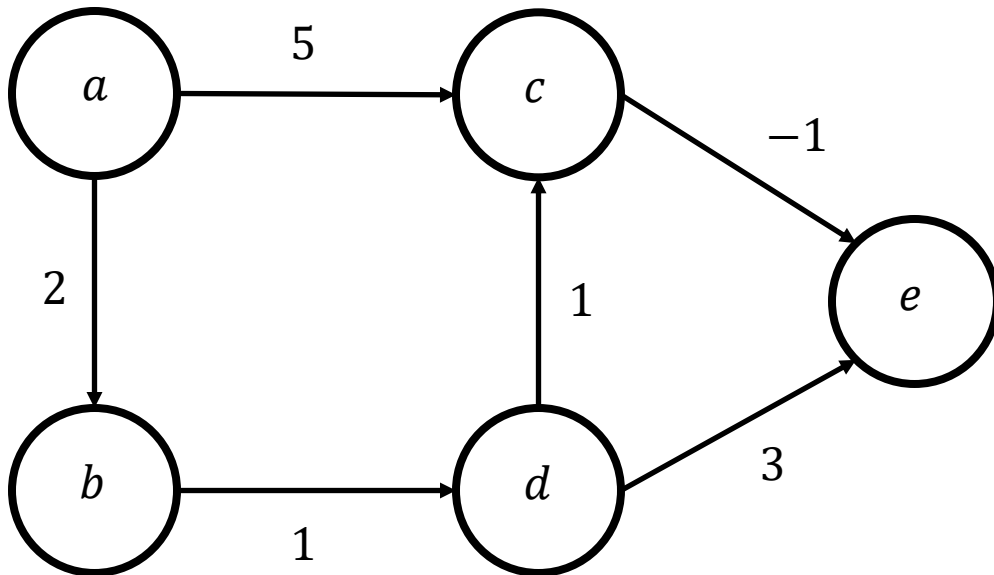
- A famous application of dynamic programming
 - Compute shortest paths from a node to all other nodes in a graph
 - E.g., all shortest paths from a to other nodes



	1 step	2 steps	
a	0	0	
b	2	2	
c	5	5	
d	∞	3	
e	∞	4	

Loop through all 1-step nodes and see if you can reach other nodes in 1 step at lower cost

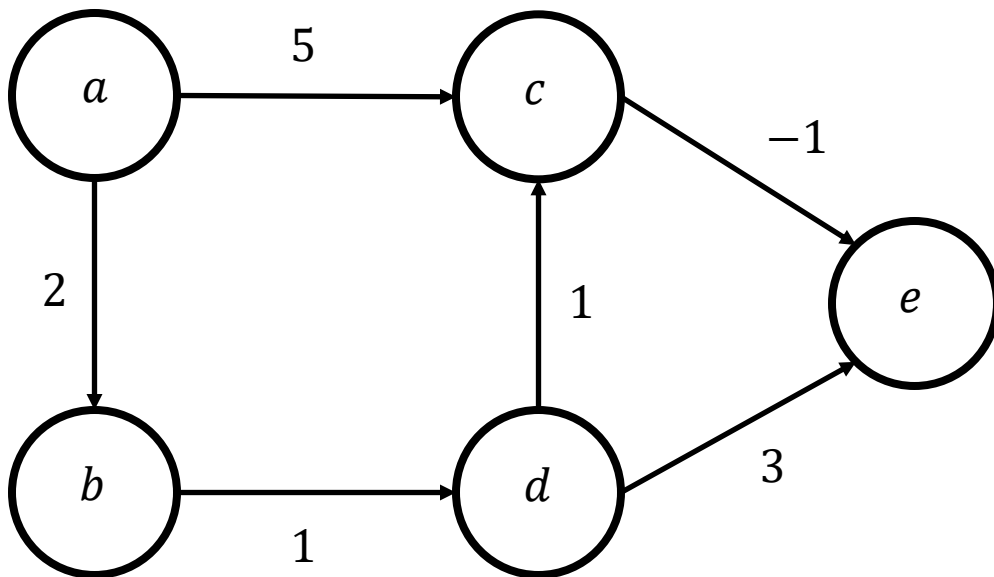
- A famous application of dynamic programming
 - Compute shortest paths from a node to all other nodes in a graph
 - E.g., all shortest paths from a to other nodes



1 step 2 steps 3 steps

a	0	0	0
b	2	2	2
c	5	5	4
d	∞	3	3
e	∞	4	4

- A famous application of dynamic programming
 - Compute shortest paths from a node to all other nodes in a graph
 - E.g., all shortest paths from a to other nodes



1 step 2 steps 3 steps 4 steps

a	0	0	0	0
b	2	2	2	2
c	5	5	4	4
d	∞	3	3	3
e	∞	4	4	3

Cost to e through c is updated to 3



- Repeat until no more improvements can be made
 - For each node n , go through all edges (n, n')
 - If $cost(a, n') > cost(a, n) + cost(n, n')$
 - Set $cost(a, n') = cost(a, n) + cost(n, n')$
- What is the worst-case complexity of the algorithm?
 - Complexity is $O(n^3)$, where n is the number of nodes
 - Each loop requires $O(n^2)$ operations
 - For each node, go through all other nodes and see if a shorter path exists
 - A total of n loops
 - Longest path to any node is n steps
- Turns out the same algorithm can be applied to MDPs
 - Find the optimal policy from any node

Workday Example

- Suppose $T = 2$

- What is $v_*^{T-1}(Pub)$?

$$q_{\pi_*}^{T-1}(Pub, Work) = -0.5$$

$$q_{\pi_*}^{T-1}(Pub, Relax) = -0.1$$

- Relax is better

- What is $v_*^{T-1}(OH)$?

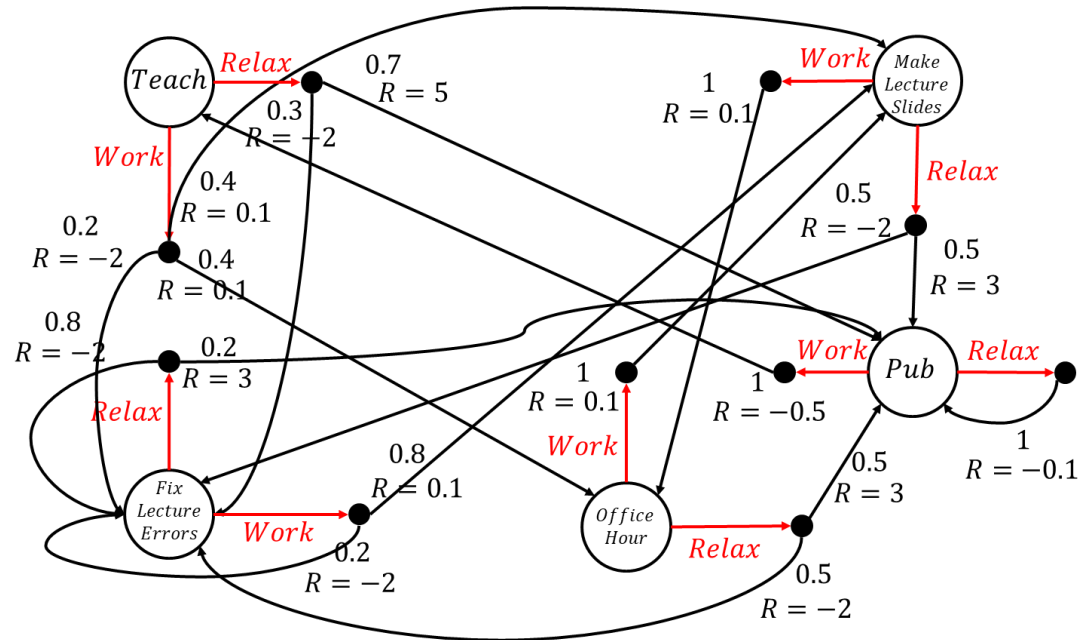
$$q_{\pi_*}^{T-1}(OH, Relax) = 0.5$$

- What is $v_*^{T-1}(FLE)$?

$$q_{\pi_*}^{T-1}(FLE, Work) = -0.32$$

- Similarly, $q_{\pi_*}^{T-1}(MLS, Relax) = 0.5$

- Similarly, $q_{\pi_*}^{T-1}(Teach, Relax) = 2.9$



Workday Example

- Suppose $\gamma = 0.9$
- What is $v_*^0(Pub)$?

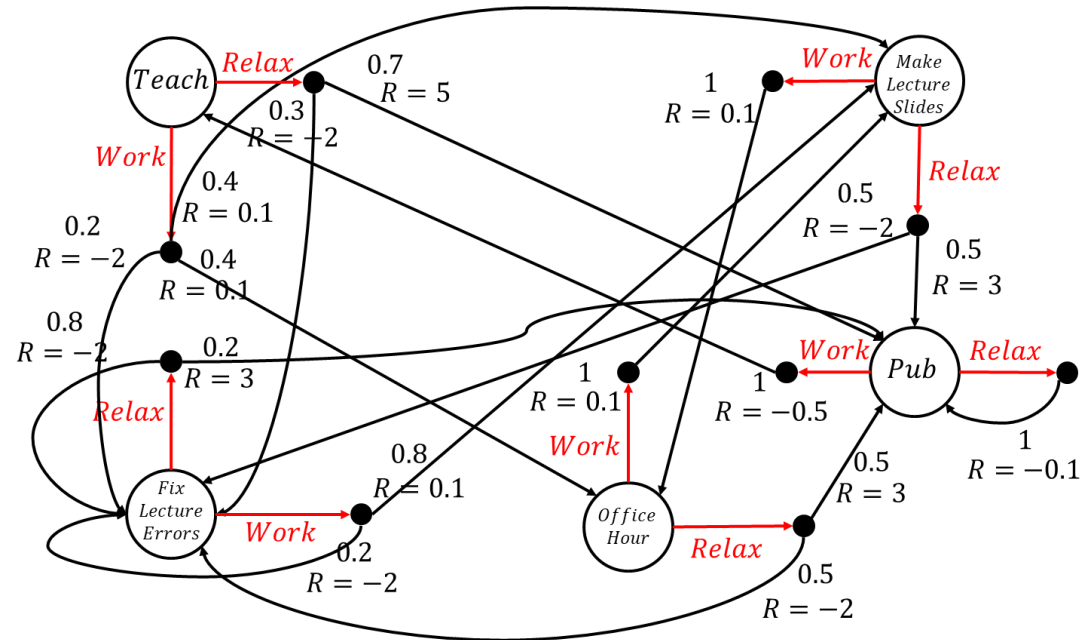
$$\begin{aligned} q_{\pi_*}^0(Pub, Work) &= \\ &= -0.5 + 0.9 * v_{\pi_*}^1(Teach) \\ &= 2.11 \end{aligned}$$

$$\begin{aligned} q_{\pi_*}^0(Pub, Relax) &= \\ &= -0.1 + 0.9 * v_{\pi_*}^1(Pub) \\ &= -0.19 \end{aligned}$$

- What is $v_*^0(FLE)$?

$$\begin{aligned} q_{\pi_*}^0(FLE, Work) &= (0.2 * -2 + 0.8 * 0.1) + \\ &= 0.2 * 0.9 * v_{\pi_*}^1(FLE) + 0.8 * 0.9 * v_{\pi_*}^1(MLS) = \\ &= -0.0176 \end{aligned}$$

$$q_{\pi_*}^0(FLE, Relax) = -1 - 0.9 * 0.276 = -1.25$$



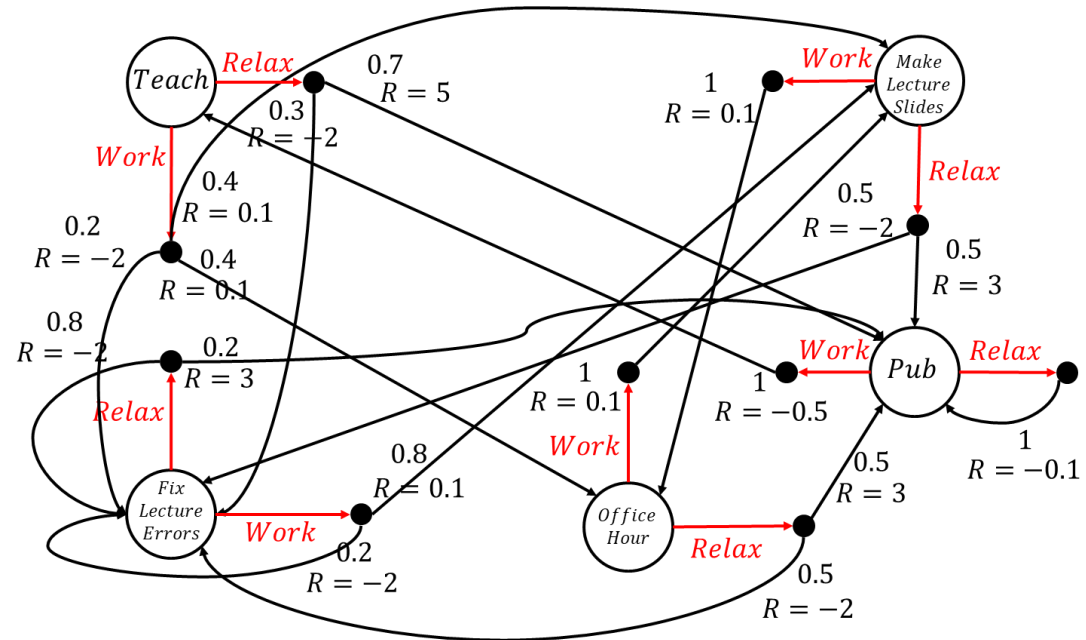
$t = 0$ $t = 1$

Teach		2.9
OH		0.5
MLS		0.5
FLE		-0.32
Pub		-0.1

Value table

Workday Example

- Suppose $\gamma = 0.9$
- The other action values computed similarly



	$t = 0$	$t = 1$	Value table
<i>Teach</i>		2.9	
<i>OH</i>		0.5	
<i>MLS</i>		0.5	
<i>FLE</i>	-0.0176	-0.32	
<i>Pub</i>	2.11	-0.1	

- Iterate backwards, starting from last decision step $T - 1$
- First, compute $q_{\pi_*}^{T-1}(s, a)$ for each state/action pair
 - i.e., compute the one-step expected reward
 - Then set $v_{\pi_*}^{T-1}(s) = \max_a q_{\pi_*}^{T-1}(s, a)$
- For $t < T - 1$, use Bellman equation to compute $q_{\pi_*}^t(s, a)$
$$q_{\pi_*}^t(s, a) = \mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}^{t+1}(S_{t+1}) | S_t = s, A_t = a]$$
 - Then set $v_{\pi_*}^t(s) = \max_a q_{\pi_*}^t(s, a)$
- What is the complexity of dynamic programming?
$$O(TS^2A)$$
 - where S is the number of states, A is the number of actions
 - for each state, loop through all actions and all other states

- So far dynamic programming only works for the case of
 - Finite horizon
 - Finite-state space
 - Finite-action space
- Finite-state and –action spaces hard to relax (for now)
- But we can modify algorithm for infinite horizon
- Policy iteration!

- Suppose we have a current policy π , potentially not optimal
 - Suppose we know $v_\pi(s), q_\pi(s, a), \forall s, a$
 - How can we improve the policy for a given s ?
 - Pick an action that has a higher q value
- We know $v_\pi(s) = q_\pi(s, \pi(s))$
 - What if there existed an action a' s.t.
$$q_\pi(s, a') \geq q_\pi(s, \pi(s))$$
 - Turns out the policy that selects a' is better
- **Policy Improvement Theorem:**
 - A policy π' is as good as, or better than, another policy π if for all $s \in S$

$$q_\pi(s, \pi'(s)) \geq v_\pi(s)$$

- First recall that for a specific action a , the q value is:

$$\begin{aligned}q_{\pi}(s, a) &= R_e(s, a) + \gamma \sum_{s'} P(s, a, s') v(s') \\ &= R_e(s, a) + \gamma \mathbf{p}(s, a)^T \mathbf{v}(s)\end{aligned}$$

– where $\mathbf{p}(s, a)^T = [P(s, a, s_1), \dots, P(s, a, s_N)]$

- Wlog, suppose π' is different from π only at s_1 , i.e.,

$$q_{\pi}(s_1, \pi'(s_1)) \geq v_{\pi}(s_1)$$

- Using the Bellman equation:

$$q_{\pi}(s_1, \pi'(s_1)) = R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T \mathbf{v}_{\pi}(s)$$

$$v_{\pi}(s_1) = q_{\pi}(s_1, \pi(s_1)) = R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T \mathbf{v}_{\pi}(s)$$

- Then

$$R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T \mathbf{v}_{\pi}(s) \geq R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T \mathbf{v}_{\pi}(s)$$

- We now construct matrix-form Bellman equations for π and π'
- For π

$$\begin{aligned} v_{\pi}(\mathbf{s}) &= \begin{bmatrix} q_{\pi}(s_1, \pi(s_1)) \\ \dots \\ q_{\pi}(s_N, \pi(s_N)) \end{bmatrix} \\ &= \begin{bmatrix} R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T v_{\pi}(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_{\pi}(\mathbf{s}) \end{bmatrix} \\ &= R_{\pi}(\mathbf{s}) + \gamma \mathbf{P}_{\pi} v_{\pi}(\mathbf{s}) \end{aligned}$$

- We now construct matrix-form Bellman equations for π and π'
- For π'

$$\begin{aligned} v_{\pi'}(\mathbf{s}) &= \begin{bmatrix} q_{\pi'}(s_1, \pi'(s_1)) \\ q_{\pi'}(s_2, \pi(s_2)) \\ \dots \\ q_{\pi'}(s_N, \pi(s_N)) \end{bmatrix} \\ &= \begin{bmatrix} R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T v_{\pi'}(\mathbf{s}) \\ R_e(s_2, \pi(s_2)) + \gamma \mathbf{p}(s_2, \pi(s_2))^T v_{\pi'}(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_{\pi'}(\mathbf{s}) \end{bmatrix} \\ &= R_{\pi'}(\mathbf{s}) + \gamma \mathbf{P}_{\pi'} v_{\pi'}(\mathbf{s}) \end{aligned}$$

– Note that \mathbf{P}_{π} and $\mathbf{P}_{\pi'}$ only differ in their first row

- Wlog, suppose π' is different from π only at s_1 , i.e.,

$$q_{\pi}(s_1, \pi'(s_1)) \geq v_{\pi}(s_1)$$

- Using the Bellman equation:

$$q_{\pi}(s_1, \pi'(s_1)) = R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T v_{\pi}(\mathbf{s})$$

$$v_{\pi}(s_1) = q_{\pi}(s_1, \pi(s_1)) = R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T v_{\pi}(\mathbf{s})$$

- Then

$$R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T v_{\pi}(\mathbf{s}) \geq R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T v_{\pi}(\mathbf{s})$$

- Stack remaining values for π in a vector as follows:

$$\begin{bmatrix} R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T v_{\pi}(\mathbf{s}) \\ R_e(s_2, \pi(s_2)) + \gamma \mathbf{p}(s_2, \pi(s_2))^T v_{\pi}(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_{\pi}(\mathbf{s}) \end{bmatrix} \geq \begin{bmatrix} R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T v_{\pi}(\mathbf{s}) \\ R_e(s_2, \pi(s_2)) + \gamma \mathbf{p}(s_2, \pi(s_2))^T v_{\pi}(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_{\pi}(\mathbf{s}) \end{bmatrix}$$

– where the inequality is interpreted element-wise

- Stack remaining values for π in a vector as follows:

$$\begin{bmatrix} R_e(s_1, \pi'(s_1)) + \gamma \mathbf{p}(s_1, \pi'(s_1))^T v_\pi(\mathbf{s}) \\ R_e(s_2, \pi(s_2)) + \gamma \mathbf{p}(s_2, \pi(s_2))^T v_\pi(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_\pi(\mathbf{s}) \end{bmatrix} \geq \begin{bmatrix} R_e(s_1, \pi(s_1)) + \gamma \mathbf{p}(s_1, \pi(s_1))^T v_\pi(\mathbf{s}) \\ R_e(s_2, \pi(s_2)) + \gamma \mathbf{p}(s_2, \pi(s_2))^T v_\pi(\mathbf{s}) \\ \dots \\ R_e(s_N, \pi(s_N)) + \gamma \mathbf{p}(s_N, \pi(s_N))^T v_\pi(\mathbf{s}) \end{bmatrix}$$

- In matrix form:

$$\begin{aligned} R_{\pi'}(\mathbf{s}) + \gamma \mathbf{P}_{\pi'} v_\pi(\mathbf{s}) &\geq R_\pi(\mathbf{s}) + \gamma \mathbf{P}_\pi v_\pi(\mathbf{s}) \\ R_{\pi'}(\mathbf{s}) + \gamma \mathbf{P}_{\pi'} v_\pi(\mathbf{s}) &\geq v_\pi(\mathbf{s}) \\ R_{\pi'}(\mathbf{s}) &\geq v_\pi(\mathbf{s}) - \gamma \mathbf{P}_{\pi'} v_\pi(\mathbf{s}) \\ R_{\pi'}(\mathbf{s}) &\geq (\mathbf{I} - \gamma \mathbf{P}_{\pi'}) v_\pi(\mathbf{s}) \end{aligned}$$

- Pre-multiply both sides by $(\mathbf{I} - \gamma \mathbf{P}_{\pi'})^{-1}$
 - Inequalities don't switch sides (don't have time to prove)

$$\begin{aligned} (\mathbf{I} - \gamma \mathbf{P}_{\pi'})^{-1} R_{\pi'}(\mathbf{s}) &\geq v_\pi(\mathbf{s}) \\ v_{\pi'}(\mathbf{s}) &\geq v_\pi(\mathbf{s}) \end{aligned}$$

Deterministic Policies: Greedy Policy Improvement

- Suppose we are given a deterministic policy π

- We can greedily improve π for each state

$$\begin{aligned}\pi'(s) &= \mathit{arg} \max_a q_\pi(s, a) \\ &= \mathit{arg} \max_a \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a]\end{aligned}$$

- By the policy improvement theorem, π' is better than or equal to π

- If $\pi' = \pi$, then $\pi' = \pi^*$:

$$v_{\pi'}(S_t) = \max_a [\mathbb{E}_\pi [R_{t+1} + \gamma v_{\pi'}(S_{t+1}) | S_t = s, A_t = a]]$$

– Bellman optimality equation!

- The greedy policy improvement approach suggests an algorithm for finding the optimal policy through iterating
 - Start from a policy, compute its value function, improve greedily for one state, repeat...

$$\pi_0 \xrightarrow{\text{E}} v_{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} v_{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \dots \xrightarrow{\text{I}} \pi_* \xrightarrow{\text{E}} v_*$$

- Terminate when you find optimal policy
- Guaranteed to terminate because there are finitely many policies in a finite-state MDP
- Policy iteration is trickier in the finite-horizon case
 - Need to evaluate the policy at each time t
 - Just use dynamic programming in the finite-horizon case

- What are the optimal actions in the long run?
 - $\pi_*(Teach) = Relax$
 - $\pi_*(OH) = Relax$
 - $\pi_*(MLS) = Relax$
 - $\pi_*(FLE) = Work$
 - $\pi_*(Pub) = Work$
- Corresponding values are
$$v_*(s) = [9.18 \ 6.31 \ 6.31 \ 5.15 \ 7.76]$$

- Start with a random policy π
- Repeat until you find the optimal policy:
 - Loop through all states
 - For each state s , loop through all actions
 - If you find an action a for which $q_{\pi}(s, a) > v_{\pi}(s)$
 - Modify π such that $\pi(s) = a$
 - Recalculate values $v_{\pi'}$ for modified policy π'
 - Go back to main loop
 - If you did not change the policy at all, terminate
 - You found the optimal!

- Policy iteration requires evaluating each new policy
 - i.e., need to compute $v_{\pi}(s)$ for all states
 - May take significant time
- Another approach is to use the Bellman optimality equation

$$\begin{aligned}v_{\pi_*}(s) &= \max_a q_{\pi_*}(s, a) \\ &= \max_a \left[\mathbb{E}_{\pi_*} [R_{t+1} + \gamma v_{\pi_*}(S_{t+1}) | S_t = s, A_t = a] \right]\end{aligned}$$

- The Bellman optimality equation suggests the recursion

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

- Starting from any v_0

- The Bellman optimality equation suggests the recursion

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

- The sequence is guaranteed to converge

- Consider the mapping

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

- Map L can be shown to be contractive

- i.e., any 2 sequences get closer to each other after each iteration

- The sequence v_k converges to a unique v^* for all v_0

- The unique v^* satisfies the Bellman optimality equation

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

- So it is the value function corresponding to the optimal policy

- Given a value function, the corresponding policy is

$$\begin{aligned}\pi(s) &= \operatorname{argmax}_a q_\pi(s, a) \\ &= \operatorname{argmax}_a \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \operatorname{argmax}_a \left[R_e(s, a) + \sum_{s'} \gamma v_\pi(s') P(s, a, s') \right]\end{aligned}$$

- Note that a v_k may not have the actual state values of the policy it represents
 - There are finitely many policies but infinitely many v_k
 - Ultimately, we don't care what the state values are as long as the policy is optimal
 - Of course, when the v_k converge, the values will converge to the values of the optimal policy

- Start from an arbitrary v_0
- For each state s , update v_{k+1} as follows:
$$v_{k+1}(s) = \max_a [\mathbb{E}[R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s, A_t = a]]$$
- Iterate until $|v_{k+1} - v_k| < \epsilon$
 - Where ϵ is a hyperparameter
 - Can use your favorite norm above, e.g., L_∞
- Unlike policy iteration, no need to invert large matrices
 - Though may require many more iterations
 - For large-state-space MDPs, value iteration likely to scale better

- The workday example has 5 states and 2 actions
 - How many policies are there in total?
$$2^5 = 32$$
 - Policy iteration likely to converge in several iterations
- Value iteration takes several dozen iterations to converge to true values
 - From an initial state of all 0's
 - Though you don't need to converge fully until you uncover the optimal policy
- Grid world has a bigger policy space
$$4^{25}$$
 - Optimal policy is very simple, so policy iteration still fast

- Dynamic programming is a powerful iterative algorithm
- Very popular in some fields of computer science and engineering
 - Widely used in control, in a similar way to RL
- Vanilla algorithm only works for finite-space MDPs
 - Overall iteration idea is still mainstream RL, however
- All algorithms discussed so far also need the user to know the MDP structure
 - Not realistic in many cases