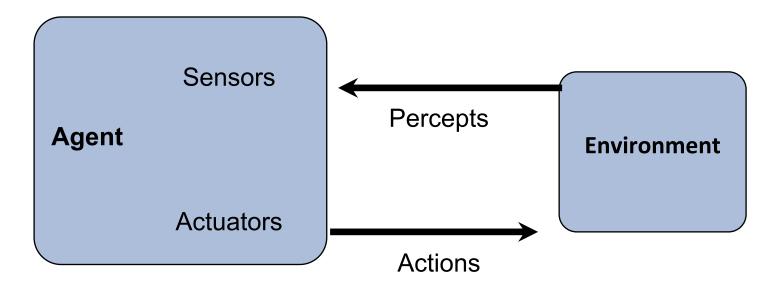


DeepMind



SO LONG CERTAINTY...



 Until now, result of taking an action in a state was deterministic



REASONING UNDER UNCERTAINTY

Learn model of outcomes

Multi-armed bandits

Reinforcement Learning

Given model of stochastic outcomes

Decision theory

Markov Decision Processes

Actions Don't Change State of the World Actions Change State of the World



EXPECTATION

 The expected value of a function of a random variable is the average, weighted by the probability distribution over outcomes

Example: expected time if take the bus

• Time: 5 min + 30 min

• Probability: 0.7 + 0.3



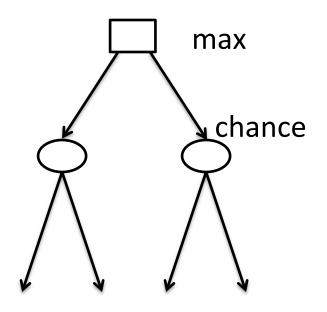
12.5 min





WHERE DO PROBABILITIES COME FROM?

- Models
- Data
- For now assume we are given the probabilities for any chance node





REASONING UNDER UNCERTAINTY

Given model of stochastic outcomes

Decision theory

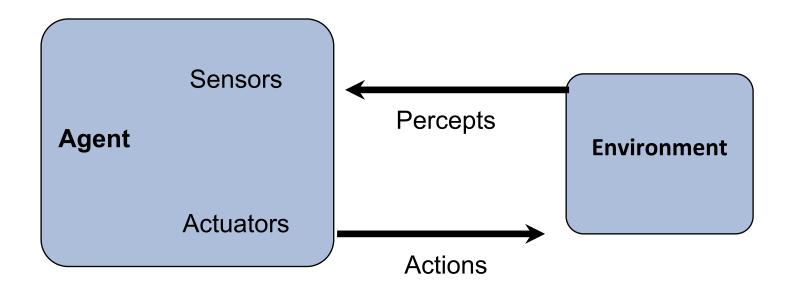
Markov Decision Processes

Actions Don't
Change State of
the World

Actions Change
State of the
World



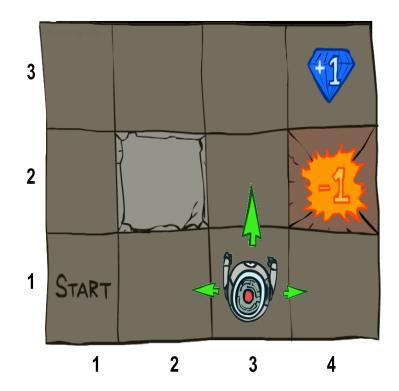
(STOCHASTICALLY) CHANGE THE WORLD



- Like planning/search, actions impact world
- But exact impact is stochastic: probability distribution over next states

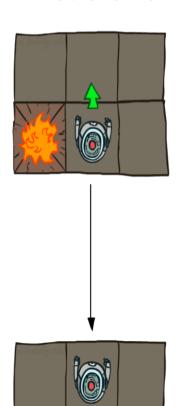
EXAMPLE: GRID WORLD

- A maze-like problem
 - The agent lives in a grid
 - Walls block the agent's path
- The agent receives rewards each time step
 - Small "living" reward each step (can be negative)
 - Big rewards come at the end (good or bad)
- Goal: maximize sum of rewards
- Noisy movement: actions do not always go as planned
 - 80% of the time, action North takes the agent North (if there is no wall there)
 - 10% of the time, North takes the agent West;
 10% East
 - If there is a wall in the direction the agent would have gone, agent stays put

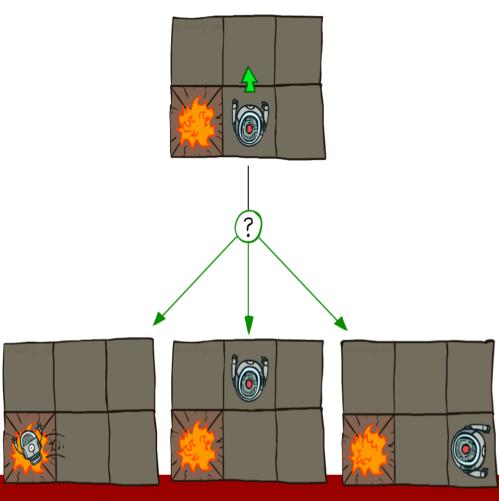


GRID WORLD ACTIONS

Deterministic Grid World



Stochastic Grid World

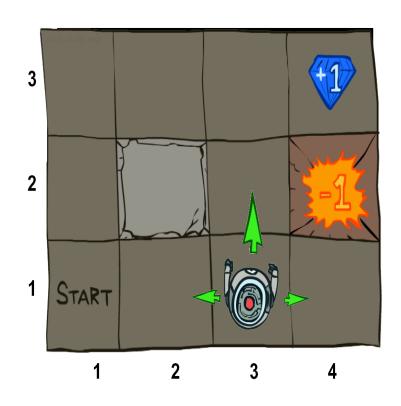


Slide adapted from Klein and Abbeel

Carnegie Mellon University

MARKOV DECISION PROCESSES

- Set of states s ∈ S
- Set of actions a ∈ A
- Transition func. T(s, a, s')
 - Probability that a from s leads to s', i.e., P(s'| s, a)
- Reward func. R(s, a, s') / R(s) / R(s,a)
- Start state or states (could be all S)
- Maybe a terminal state
- Discount factor
- MDPs are non-deterministic search problems





MARKOV DECISION PROCESSES







MARKOV PROPERTY

- Called Markov decision process because the outcome of an action depends only on the current state
- $p(s_{t+1}|s_1,a_1,s_2,a_2,...s_t,a_t)=p(s_{t+1}|s_t,a_t)$

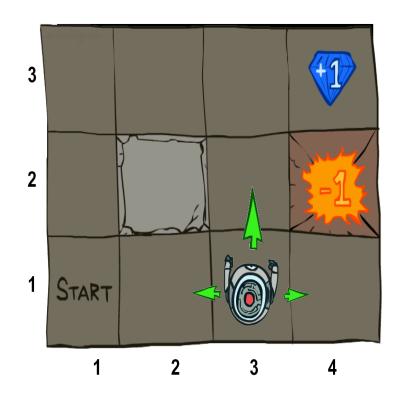


POLICIES

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal
- In MDPs instead of plans, we have a policies
- A policy π*: S → A
 - Specifies what action to take in each state

POLL: HOW MANY POLICIES?

- How many non-terminal states?
- How many actions?
- How many deterministic policies over non-terminal states?





OPTIMAL POLICIES

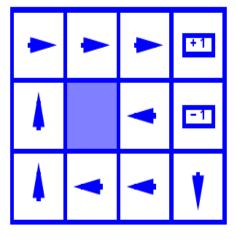
- Optimal plan had minimal cost to reach goal
- Utility or value of a policy π starting in state s is the expected sum of future rewards will receive by following π starting in state s

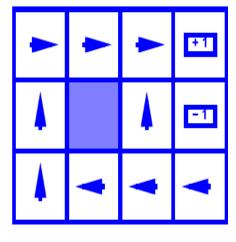
 Optimal policy has maximal expected sum of rewards from following it



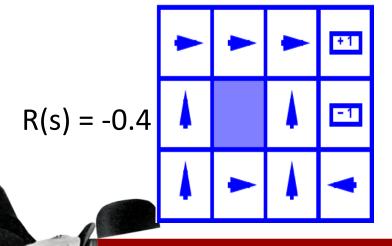
OPTIMAL POLICIES

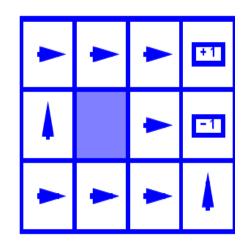






$$R(s) = -0.03$$

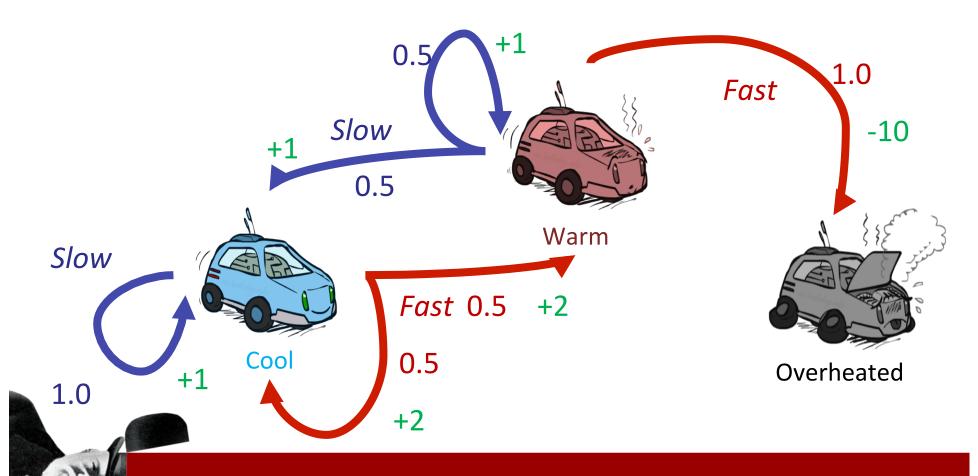




$$R(s) = -2.0$$

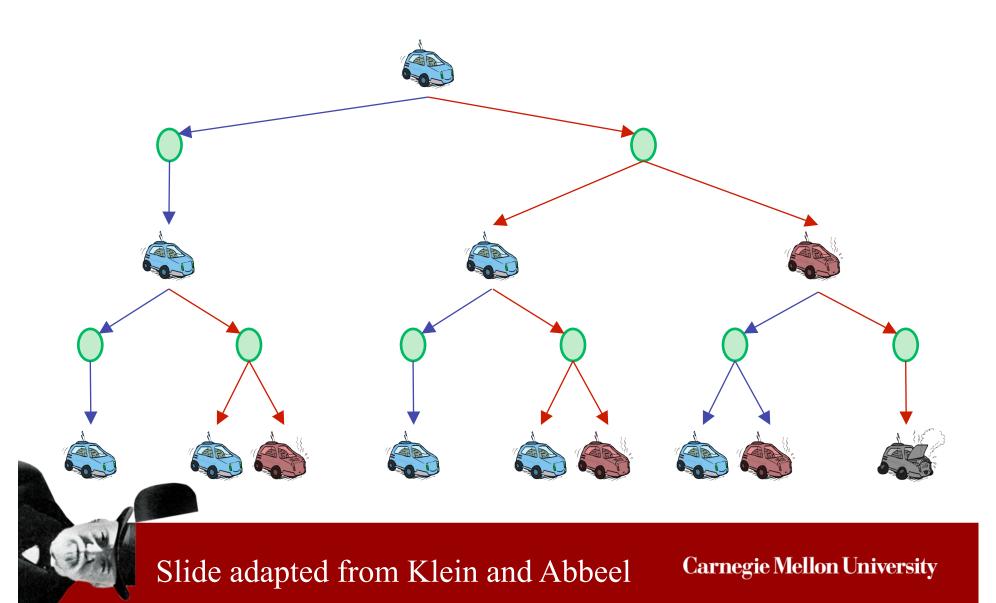
- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: Slow, Fast
- Going faster gets double reward

Example: Racing

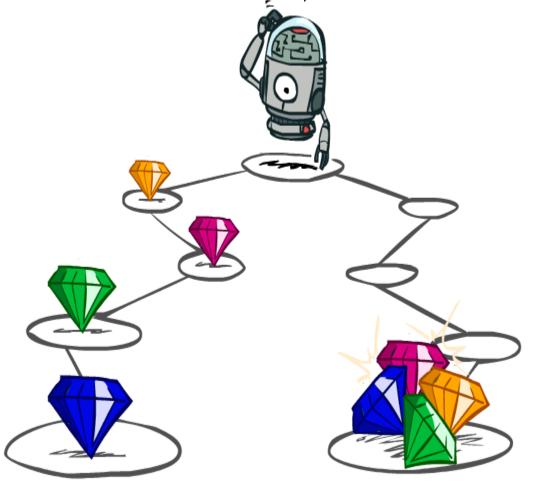


Slide adapted from Klein and Abbeel

RACING SEARCH TREE



UTILITIES OF SEQUENCES



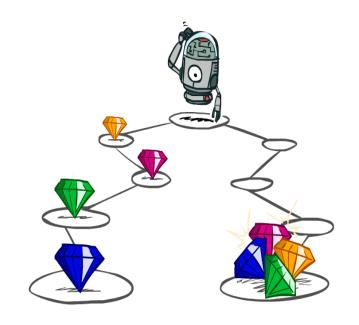
UTILITIES OF SEQUENCES

 What preferences should an agent have over reward sequences?

More or less?

Now or later?

[0, 0, 1] or [1, 0, 0]





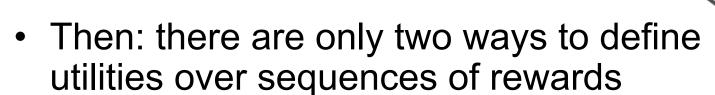
STATIONARY PREFERENCES

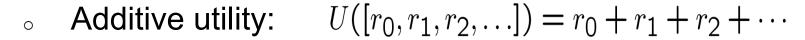
Theorem: if we assume stationary preferences:

$$[a_1, a_2, \ldots] \succ [b_1, b_2, \ldots]$$

$$\updownarrow$$

$$[r, a_1, a_2, \ldots] \succ [r, b_1, b_2, \ldots]$$





Discounted utility: $U([r_0, r_1, r_2, ...]) = r_0 + \gamma r_1 + \gamma^2 r_2 ...$

WHAT ARE DISCOUNTS?

- It's reasonable to prefer rewards now to rewards later
- Decay rewards exponentially



Worth _Now Worth Next Step Worth In Two Steps

DISCOUNTING $U([r_0, r_1, r_2, ...]) = r_0 + \gamma r_1 + \gamma^2 r_2 ...$

Given:

10				1
а	b	С	d	е

- Actions: East, West
- Terminal states: a and e (end when reach one or the other)
- Transitions: deterministic
- Reward for reaching a is 10 (regardless of initial state & action, e.g. r(s,action,a) = 10), reward for reaching e is 1, and the reward for reaching all other states is 0
- Quiz 1: For γ = 1, what is the optimal policy?
- Quiz 2: For γ = 0.1, what is the optimal policy for states b, c and d?
- Quiz 3: For which γ are West and East equally good when in state d?



INFINITE UTILITIES?!

- Problem: What if the game lasts forever? Do we get infinite rewards?
- Solutions:
 - Finite horizon: (similar to depth-limited search)
 - Terminate episodes after a fixed T steps (e.g. li.,
 - Gives nonstationary policies (π depends on time left)
 - Discounting: use $0 < \gamma < 1$

$$U([r_0, \dots r_\infty]) = \sum_{t=0}^{\infty} \gamma^t r_t \le R_{\text{max}}/(1-\gamma)$$

- Smaller γ means smaller "horizon" shorter term focus
- Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "overheated" for racing)

RECAP: DEFINING MDPs

- Markov decision processes:
 - Set of states S
 - Start state s₀
 - Set of actions A
 - Transitions P(s'|s,a) (or T(s,a,s'))
 - Rewards R(s,a,s') (and discount γ)
- MDP quantities so far:
 - Policy = Choice* of action for each state
 - Utility/Value = sum of (discounted) rewards



VALUE OF A POLICY IN EACH STATE

- Expected immediate reward for taking action prescribed by policy π for that state
- And expected future reward get after taking that action from that state and following π

$$V^{\pi}(s) = \sum_{s' \in S} p(s' \mid s, \pi(s)) \Big[R(s, \pi(s), s') + \gamma V^{\pi}(s') \Big]$$

 Future reward depends on horizon (how many more steps get to act). For now assume infinite

Q: STATE-ACTION VALUE

- Expected immediate reward for taking action
- And expected future reward get after taking that action from that state and following π

$$Q^{\pi}(s,a) = \sum_{s' \in S} p(s' \mid s,a) \left[R(s,a,s') + \gamma V^{\pi}(s') \right]$$



OPTIMAL VALUE V* AND π^*

- Optimal value: Highest possible value for each s
- Satisfies the Bellman Equation

$$V^*(s_i) = \max_{a} \left(\sum_{s_j \in S} p(s_j \mid s_i, a) \left[R(s_i, a, s') + \gamma V^*(s_j) \right] \right)$$

Optimal policy

$$\pi^*(s_i) = \operatorname*{argmax}_{a} Q(s_i, a)$$

$$= \underset{a}{\operatorname{argmax}} \left(\sum_{s_j \in S} p(s_j \mid s_i, a) \left[R(s_i, a, s') + \gamma V * (s_j) \right] \right)$$

Want to find these optimal values!

VALUE ITERATION

Bellman equation inspires an update rule

$$V * (s_i) = \max_{a} \left(\sum_{s_j \in S} p(s_j \mid s_i, a) \left[R(s, a, s') + \gamma V * (s_j) \right] \right)$$



$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) \left[R(s, a, s') + \gamma V_{k-1}(s_{j}) \right] \right)$$

Form of dynamic programming



ALSO CALLED A BELLMAN BACKUP

$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) \left[R(s, a, s') + \gamma V_{k-1}(s_{j}) \right] \right)$$

 In shorthand, for performing the above computation for all states,

$$V_k = BV_{k-1}$$



VALUE ITERATION ALGORITHM

- 1. Initialize V₀(s_i)=0 for all states s_{i.} Set K=1
- 2. While k < desired horizon or (if infinite horizon) values have converged
 - For all s,

$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) \left[R(s, a, s') + \gamma V_{k-1}(s_{j}) \right] \right)$$

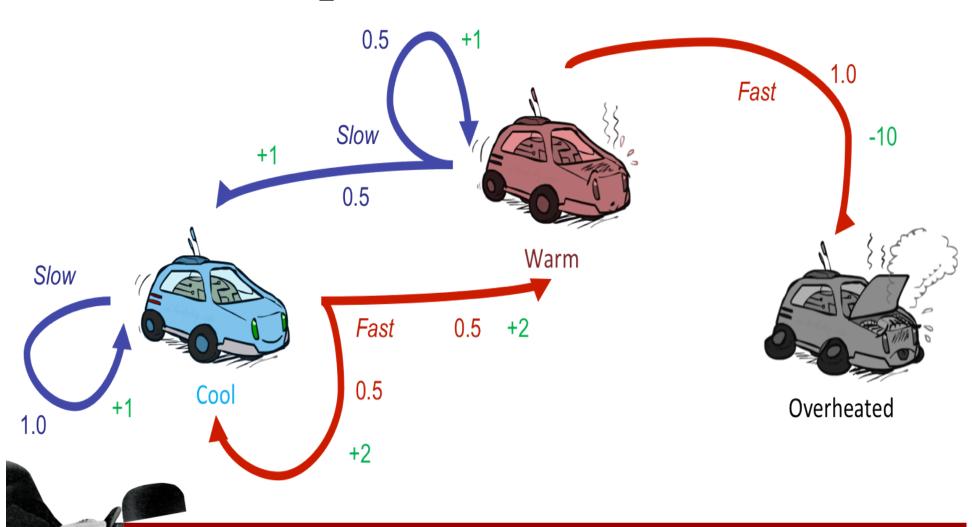
3. Extract Policy

$$\pi_k(s_i) = \underset{a}{\operatorname{argmax}} \left(\sum_{s_j \in S} p(s_j \mid s_i, a) \left[R(s, a, s') + \gamma V_{k-1}(s_j) \right] \right)$$

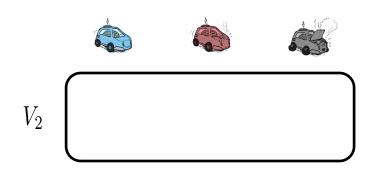


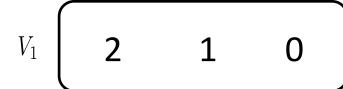
Calculate V₂(warmCar)

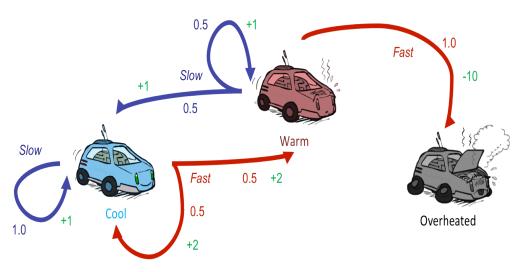
Assume Y=1



For General Practice, check can calculate V₂(warmCar)







Assume Y=1

$$V_0$$
 \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc \bigcirc

$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) \left[R(s, a, s') + \gamma V_{k-1}(s_{j}) \right] \right)$$

COMPUTATIONAL COST FOR 1 UPDATE OF V(S) FOR ALL S IN VALUE ITERATION?

For all s,

$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) \left[R(s, a, s') + \gamma V_{k-1}(s_{j}) \right] \right)$$



WILL VALUE ITERATION CONVERGE FOR INFINITE HORIZON PROBLEMS?



CONTRACTION OPERATOR

- Let O be an operator
- If |OV OV'| <= |V-V|'
- Then O is a contraction operator



WILL VALUE ITERATION CONVERGE?

- Yes, if discount factor γ < 1 or end up in a terminal state with probability 1
- Bellman equation is a contraction if discount factor, γ < 1
- If apply it to two different value functions, distance between value functions shrinks after apply Bellman equation to each



PROPERTIES OF CONTRACTION

- Only has 1 fixed point (the point reach if apply a contraction operator many times)
 - If had two, then would not get closer when apply contraction function, violating definition of contraction
- When apply contraction function to any argument, value must get closer to fixed point
 - Fixed point doesn't move
 - Repeated function applications yield fixed point



VI CONVERGES

 Value iteration converges to unique solution which is optimal value function

$$\lim_{k\to\infty} V_k = V^*$$

$$\begin{aligned} \|V_{k+1} - V^*\|_{\infty} &= \|BV_k - V^*\|_{\infty} \le \gamma \|V_k - V^*\|_{\infty} \le \dots \\ &\le \gamma^{k+1} \|V_0 - V^*\|_{\infty} \to 0 \end{aligned}$$



DISCUSS AND REPORT BACK: DOES INITIALIZATION IMPACT FINAL

VALUE?

Value Iteration Algorithm

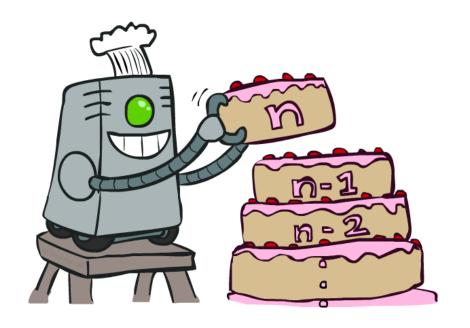
- Init $V_0(s_i)$ for all states s_i
- 2. k=1
- 3. While k < desired horizon or (if infinite horizon) values have converged
 - For all s,

$$V_{k}(s_{i}) = \max_{a} \left(\sum_{s_{j} \in S} p(s_{j} \mid s_{i}, a) [R(s, a, s') + \gamma V_{k-1}(s_{j})] \right)$$



VALUE ITERATION IN INFINITE HORIZON

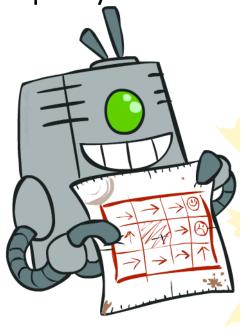
- Always have optimal values as if had only t more decisions to make
- Extracting optimal policy for t-th step yields optimal action should take, if have t more steps to act
 - But not for (t-1), (t-2)... steps
- Before convergence, these are approximations (because actually get to act forever!)
 - After convergence, value is always the same if do another update, and so is the policy







Policy Iteration
Maintain value of
following a particular
policy forever

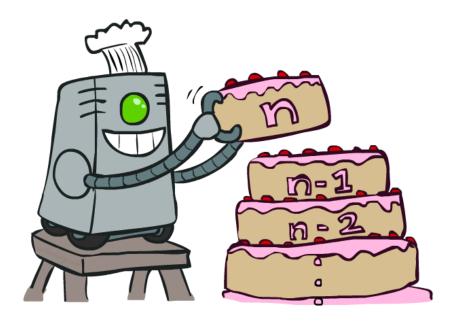


Value Iteration

Maintain optimal

values if have n more

actions

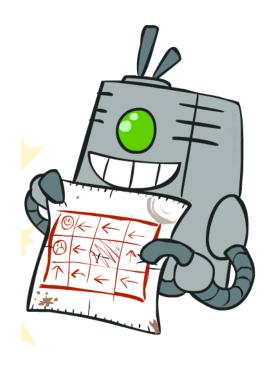


Drawings by Ketrina Yim



POLICY ITERATION FOR INFINITE HORIZON

- Instead of maintaining optimal value if have t steps left...
- Calculate exact value of acting in infinite horizon for a particular policy
- 2. Then try to improve the policy
- 3. Repeat 1 & 2 until policy doesn't change

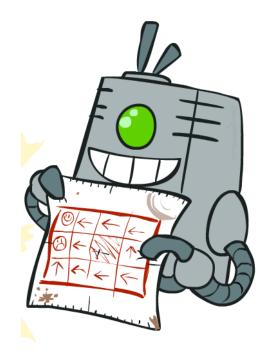


Drawing by Ketrina Yim



POLICY ITERATION FOR INFINITE HORIZON

- Instead of maintaining optimal value if have t steps left...
- Calculate exact value of acting in infinite horizon for a particular policy (evaluation)
- 2. Then try to improve the policy (improvement)
- 3. Repeat 1 & 2 until policy doesn't change



Drawing by Ketrina Yim



REVIEW: VALUE OF A POLICY

- Expected immediate reward for taking action prescribed by policy
- And expected future reward get after taking policy from that state and following π

$$V^{\pi}(s) = \sum_{s' \in S} p(s' \mid s, \pi(s)) \Big[R(s, \pi(s), s') + \gamma V^{\pi}(s') \Big]$$



POLICY EVALUATION

Goal: compute V^π(s) for all s



IDEA 1: USE SMALL VARIANT OF VALUE ITERATION

- Initialize V₀(s) to 0 for all s
- For k=1... convergence

$$V_{k}^{\pi}(s_{i}) = \sum_{s_{i} \in S} p(s_{j} \mid s_{i}, \pi(s_{i})) \left[R(s, \pi(s), s') + \gamma V_{k-1}^{\pi}(s_{j}) \right]$$



IDEA 2: NO MAX IN EQN SO LINEAR SET OF EQUATIONS... ANALYTIC SOLUTION!

$$V^{\pi}(s) = \sum_{s' \in S} p(s' \mid s, \pi(s)) \Big[R(s, \pi(s), s') + \gamma V^{\pi}(s') \Big]$$

Let T^{π} be a S x S matrix where the (i,j) entry is:

$$T^{\pi}(s_i, s_j) = p(s_j \mid s_i, \pi(s_i))$$

$$\vec{V} = T^{\pi} \vec{R} + \gamma T^{\pi} \vec{V}$$

$$\vec{V} - \gamma T^{\pi} \vec{V} = T^{\pi} \vec{R}$$

$$\vec{V} = (1 - \gamma T^{\pi})^{-1} T^{\pi} \vec{R}$$

Requires taking an inverse of a S by S matrix O(S³)



POLICY IMPROVEMENT

- Have V^π(s) for all s
- Want a better policy
- Idea:
 - Find the state-action Q value of doing an action followed by following π forever, for each state
 - Then take argmax of Qs



POLICY IMPROVEMENT

- Have V^π(s) for all s
- First compute

$$Q^{\pi}(s,a) = \sum_{s' \in S} p(s' \mid s,a) \left[R(s,a,s') + \gamma V^{\pi}(s') \right]$$

Then extract new policy.
 For each s,

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



DELVING DEEPER INTO IMPROVEMENT STEP

$$Q^{\pi}(s,a) = \sum_{s' \in S} p(s' \mid s,a) \left[R(s,a,s') + \gamma V^{\pi}(s') \right]$$

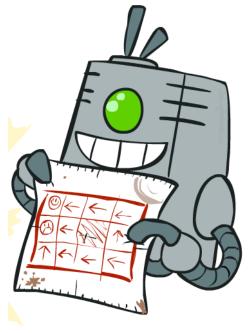
$$\max_{a} Q^{\pi}(s, a) \ge V^{\pi}(s)$$

$$\pi'(s) = \operatorname{arg\,max}_a Q^{\pi}(s, a)$$

- So if, starting at any state, we followed π '(s), transitioned to a new state s', and then followed π forever, our expected sum of rewards would be at least as good as if we had always followed π
- But we're not doing that, we're getting a new policy and proposing always following that policy π'...

POLICY ITERATION FOR INFINITE HORIZON

- 1. Policy Evaluation: Calculate exact value of acting in infinite horizon for a particular policy
- 2. Policy Improvement
- 3. Repeat 1 & 2 until policy doesn't change







DISCUSS AND REPORT BACK: IF POLICY DOESN'T CHANGE (π'(s) =π(s) FOR ALL s), CAN IT EVER CHANGE AGAIN IN MORE ITERATIONS?

Given V^π(s) for all s

$$Q^{\pi}(s,a) = \sum_{s' \in S} p(s' \mid s,a) \Big[R(s,a,s') + \gamma V^{\pi}(s') \Big]$$

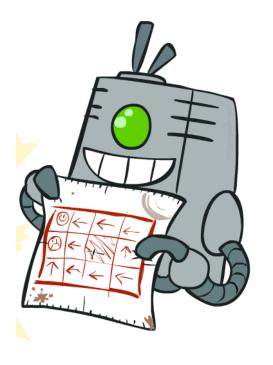
$$\pi'(s) = \arg\max_{a} Q^{\pi}(s, a)$$



POLICY ITERATION FOR

INFINITE HORIZON

- 1. Policy Evaluation:
 - Calculate exact value of acting in infinite horizon for a particular policy
- 2. Policy Improvement
- 3. Repeat 1 & 2 until policy doesn't change



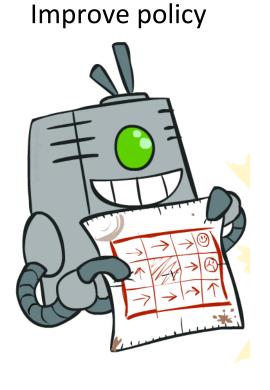


Drawing by Ketrina Yim

DISCUSS AND REPORT BACK POLICY ITERATION IS GUARANTEED TO CONVERGE. THINK ABOUT WHY, AND HOW MANY ITERATIONS COULD IT TAKE (HOW MANY POLICIES **COULD WE HAVE TO SORT** THROUGH)?

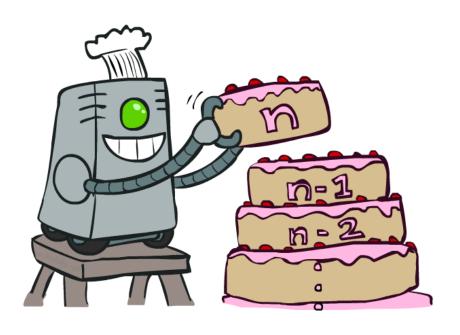


Policy Iteration Maintain value of policy



Value Iteration

Keep optimal value for finite steps, increase steps



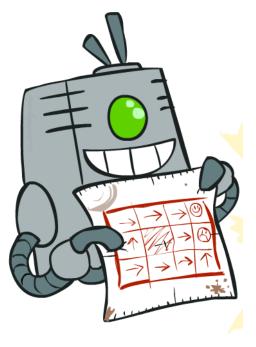
Drawings by Ketrina Yim



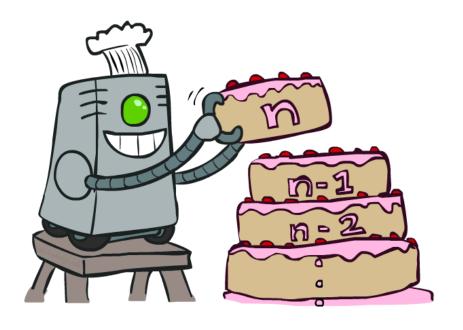
Policy Iteration

Fewer Iterations

More expensive per iteration



Value Iteration More iterations Cheaper per iteration



Drawings by Ketrina Yim



MDPs: What You Should Know

- Definition
- How to define for a problem
- Value iteration and policy iteration
 - How to implement
 - Convergence guarantees
 - Computational complexity

