



CMU

MDPs

15-381 / 781

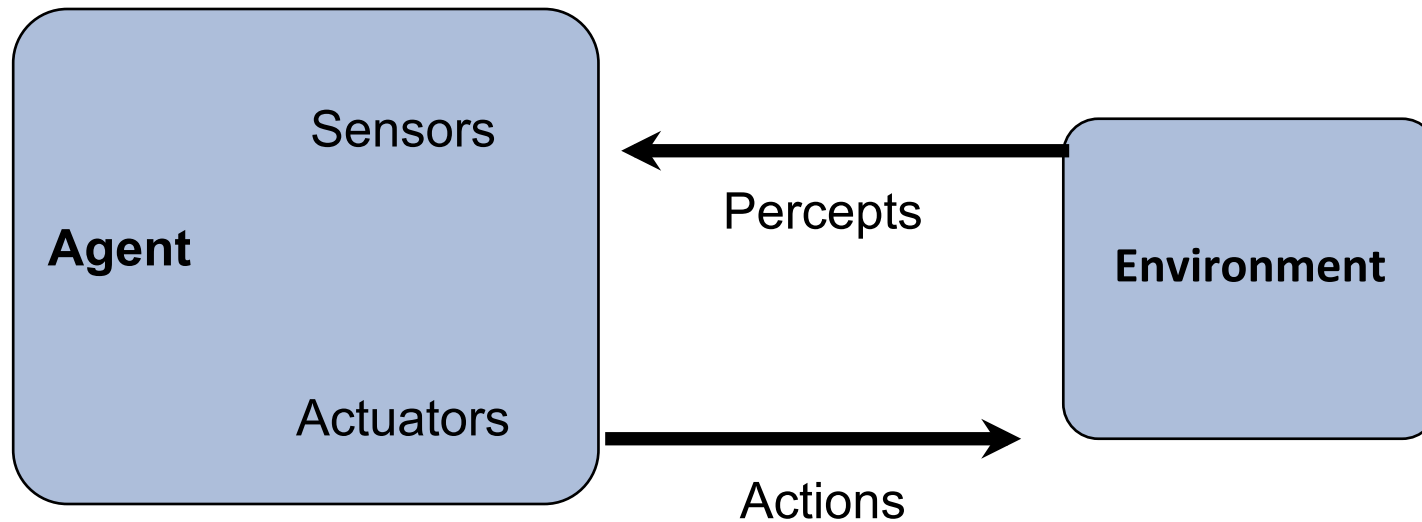
EMMA BRUNSKILL (THIS TIME)

ARIEL PROCACCIA

- 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

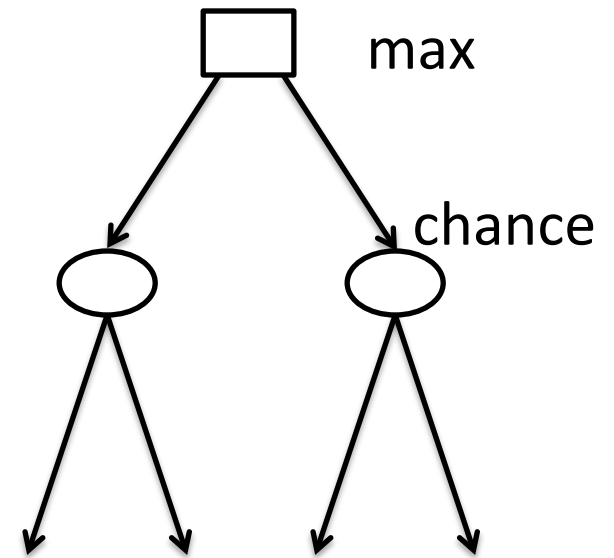
- Example: expected time if take the bus

- Time: 5 min + 30 min
- Probability: 0.7 + 0.3 \Rightarrow 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

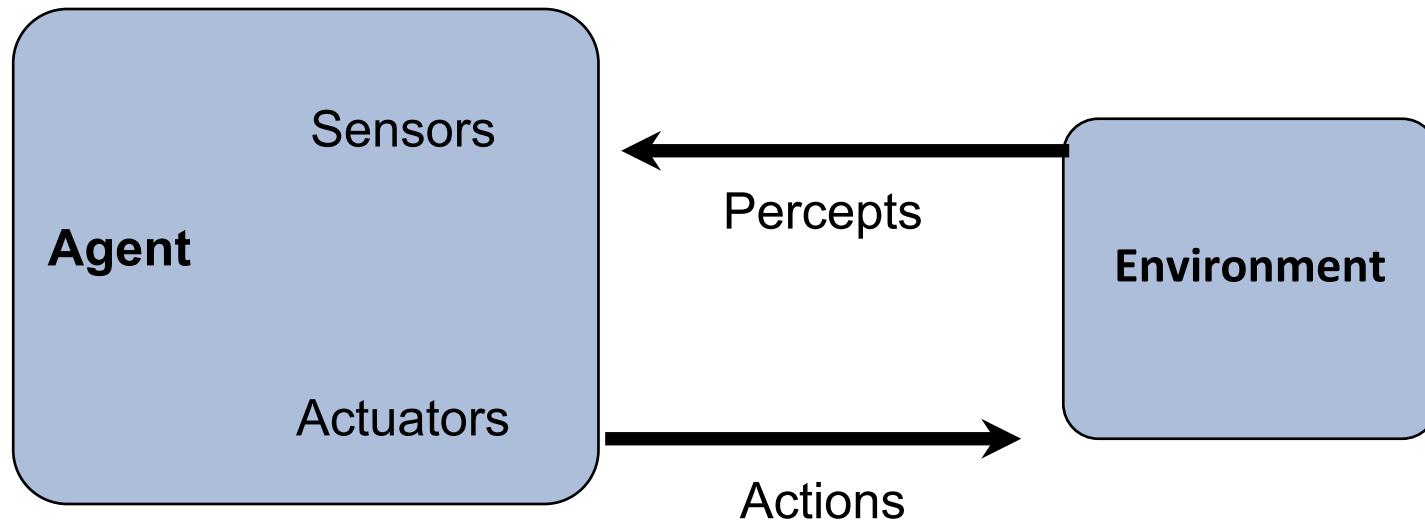
Decision theory	Markov Decision Processes

Actions Don't
Change State of
the World

Actions Change
State of the
World



(STOCHASTICALLY) CHANGE THE WORLD

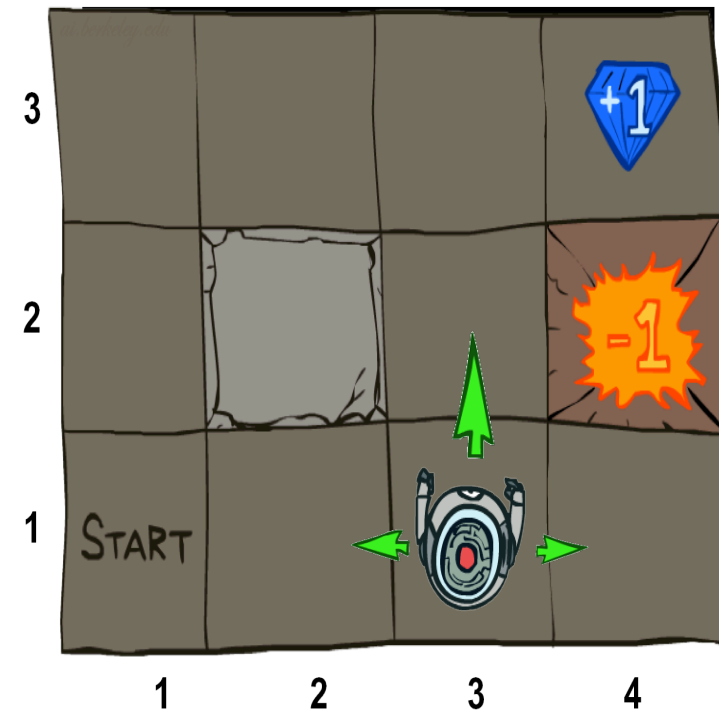


- Like planning/search, actions impact world
- But exact impact is stochastic: probability distribution over next states



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

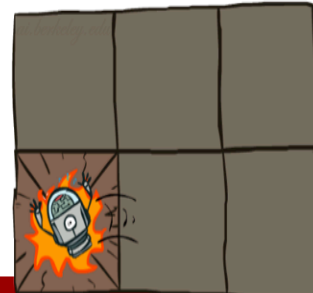
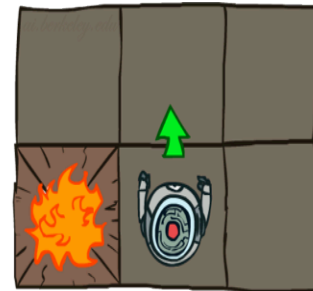


GRID WORLD ACTIONS

Deterministic Grid World



Stochastic Grid World

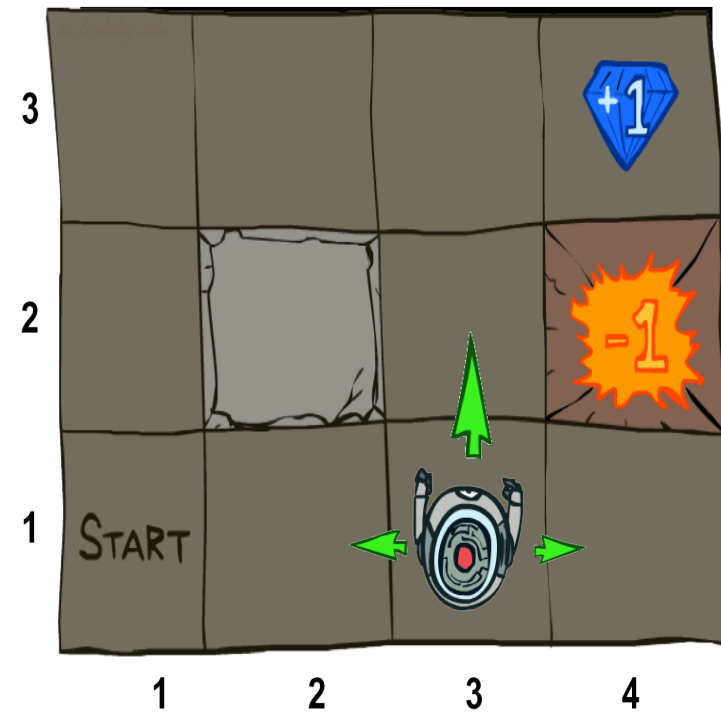


Slide adapted from Klein and Abbeel

Carnegie Mellon University

MARKOV DECISION PROCESSES

- Set of states $s \in S$
- Set of actions $a \in A$
- 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
- $p(s_{t+1}|s_1, a_1, s_2, a_2, \dots, s_t, a_t) = p(s_{t+1}|s_t, a_t)$



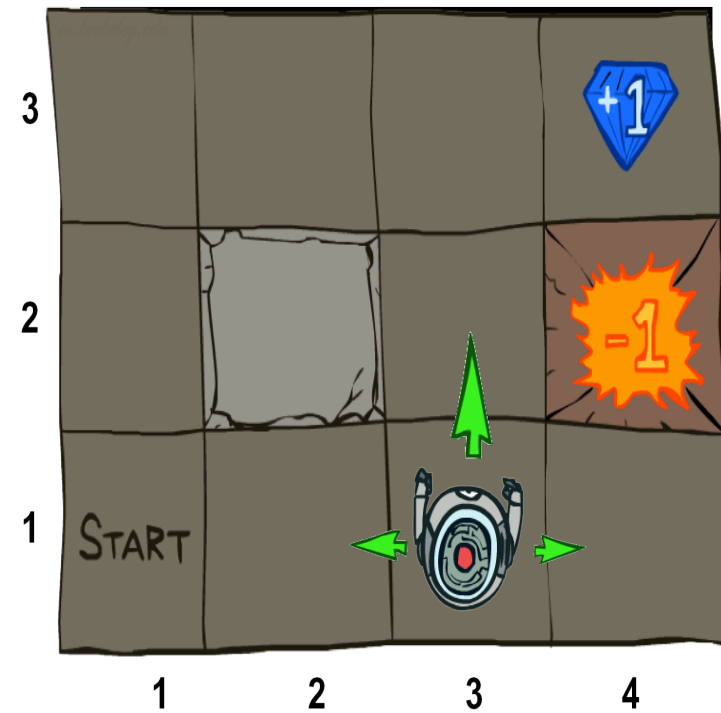
POLICIES

- 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 $\pi^*: S \rightarrow A$
 - Specifies what action to take in each state



HOW MANY POLICIES?

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



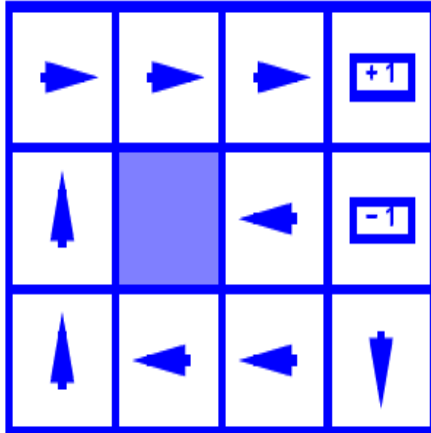
OPTIMAL POLICIES

- 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

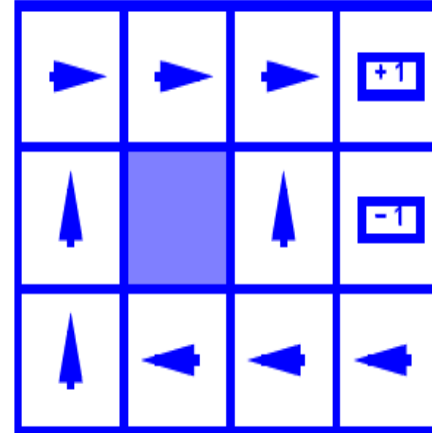


OPTIMAL POLICIES

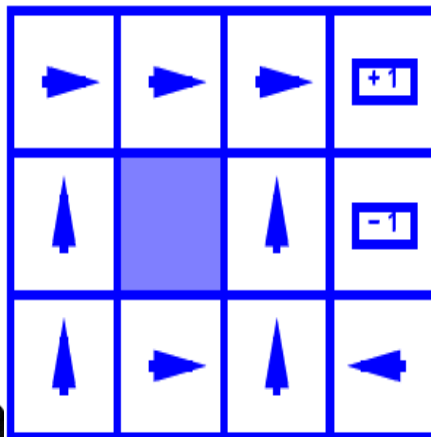
$R(s) =$
-0.01



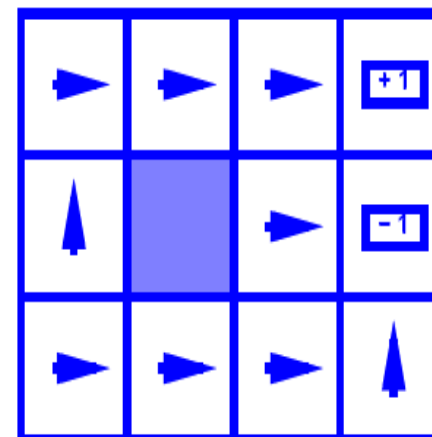
$R(s) =$
-0.03



$R(s) =$ -0.4



$R(s) =$ -2.0

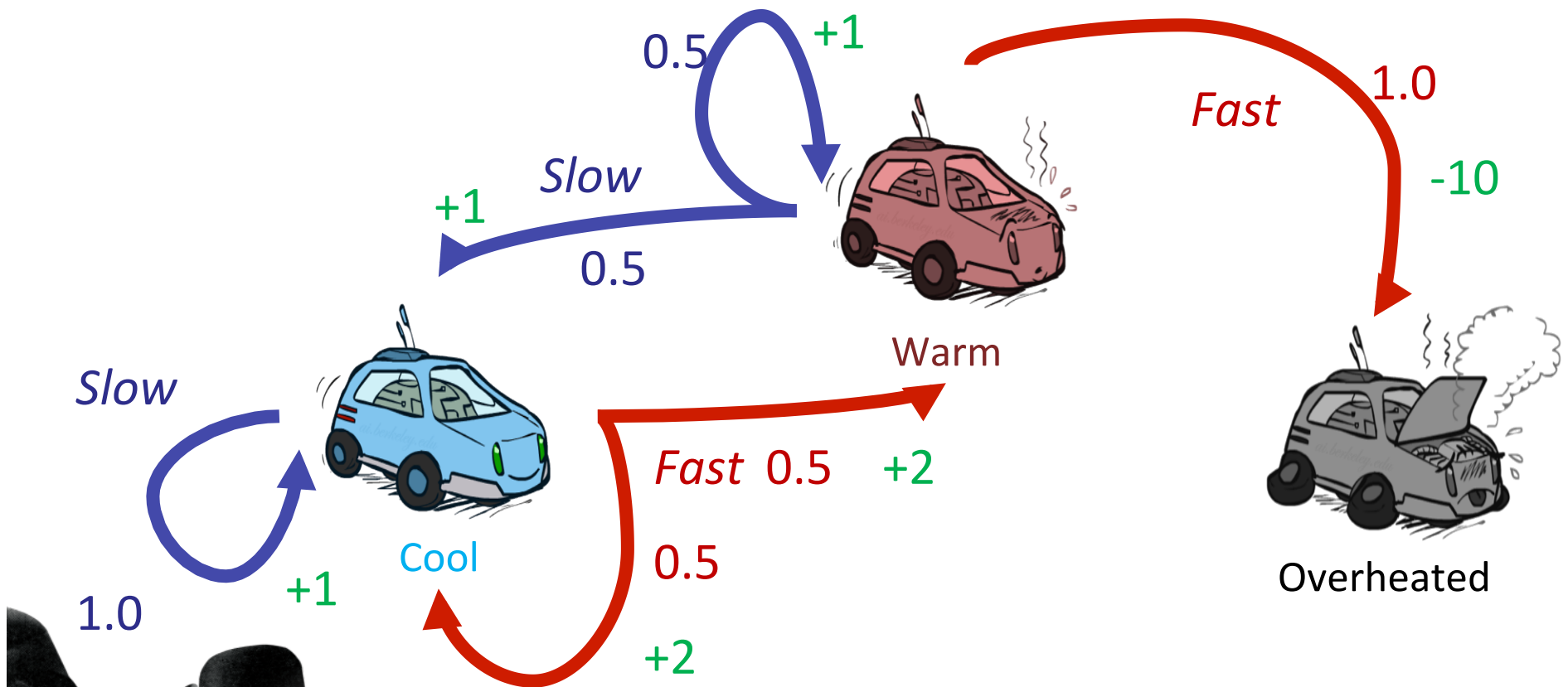


Slide adapted from Klein and Abbeel

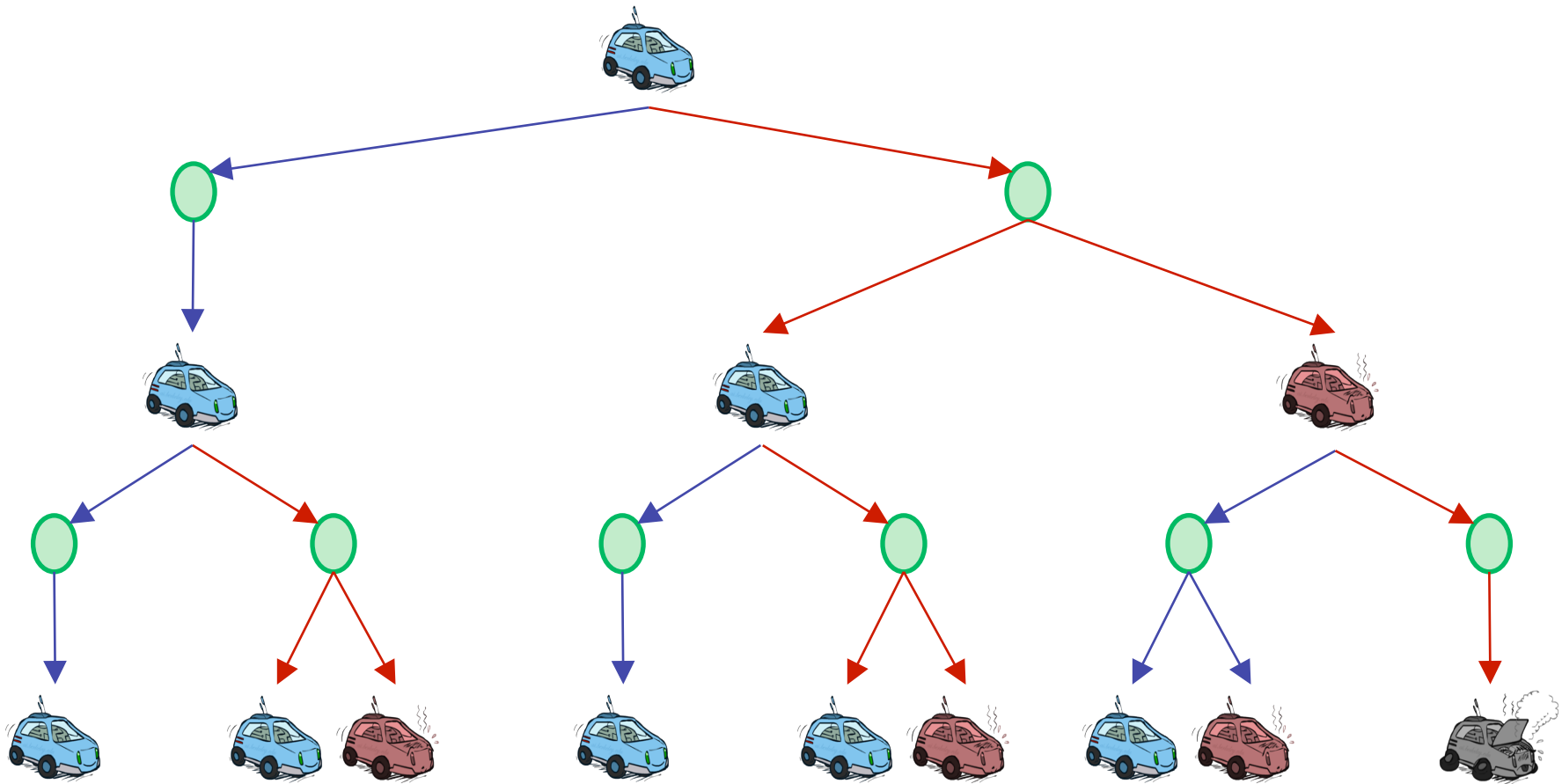
Carnegie Mellon University

Example: Racing

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



RACING SEARCH TREE



Slide adapted from Klein and Abbeel

Carnegie Mellon University

UTILITIES OF SEQUENCES



Slide adapted from Klein and Abbeel

Carnegie Mellon University

UTILITIES OF SEQUENCES

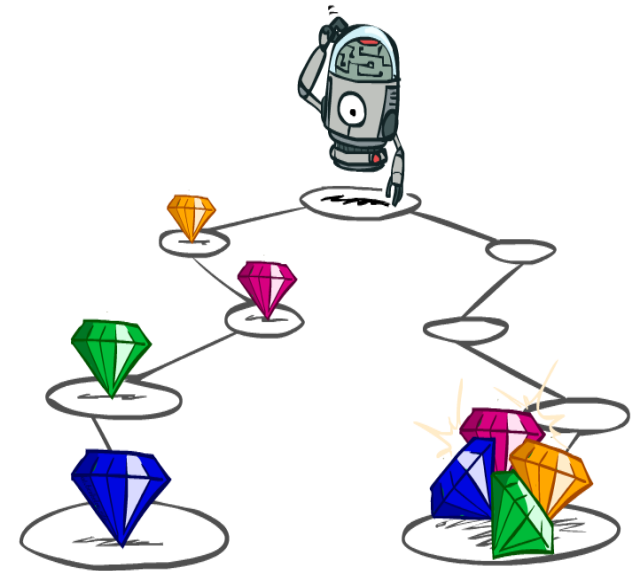
- What preferences should an agent have over reward sequences?

- More or less?

[1, 2, 2] or [2, 3, 4]

- Now or later?

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



STATIONARY PREFERENCES

- Theorem: if we assume stationary preferences:

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

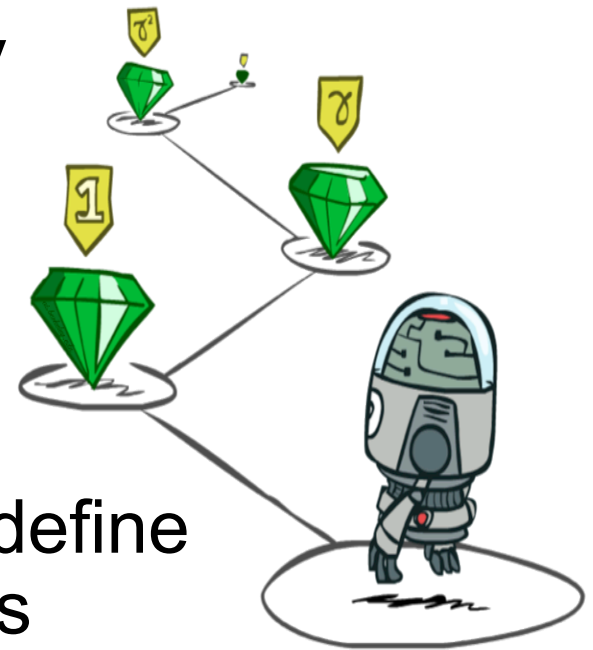


$$[r, a_1, a_2, \dots] \succ [r, b_1, b_2, \dots]$$

- Then: there are only two ways to define utilities over sequences of rewards

- Additive utility: $U([r_0, r_1, r_2, \dots]) = r_0 + r_1 + r_2 + \dots$

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



WHAT ARE DISCOUNTS?

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



1

Worth
Now



γ

Worth Next
Step



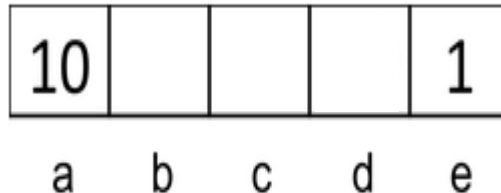
γ^2

Worth In Two
Steps



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

- Given:

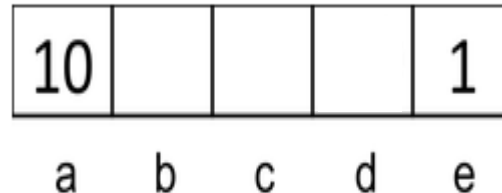


- Actions: East, West
 - Terminal states: a and e (end when reach one or the other)
 - Transitions: deterministic
 - Reward for reaching a is 10 (regardless of initial state & action, e.g. $r(s, \text{action}, a) = 10$), reward for reaching e is 1, and the reward for reaching all other states is 0
- Quiz 1: For $\gamma = 1$, what is the optimal policy?
 - Quiz 2: For $\gamma = 0.1$, 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?



QUIZ: DISCOUNTING

- Given:



- 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 a & action, e.g. $r(s, \text{action}, a) = 10$), reward for reaching e is 1, and the reward for reaching all other states is 0
- Quiz 1: For $\gamma = 1$, what is the optimal policy?
 - In all states, Go West (towards a)
- Quiz 2: For $\gamma = 0.1$, what is the optimal policy?
 - $b=W, c=W, d=E$
- Quiz 3: For which γ are West and East equally good when in state d? $\text{Gamma} = \sqrt{1/10}$



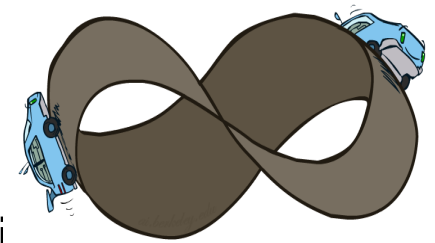
INFINITE UTILITIES?!

- 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 $0 < \gamma < 1$

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

- 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 s_0
 - 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

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

- 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

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



OPTIMAL VALUE V^* AND π^*

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

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

- Optimal policy

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

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

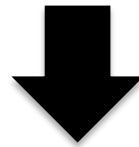
- Want to find these optimal values!



VALUE ITERATION

- Bellman equation inspires an update rule

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



$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$

- Form of *dynamic programming*



ALSO CALLED A BELLMAN BACKUP

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$

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

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



VALUE ITERATION ALGORITHM

1. Initialize $V_0(s_i)=0$ for all states s_i , Set $K=1$
2. While $k <$ desired horizon or (if infinite horizon) values have converged

- For all s ,

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$

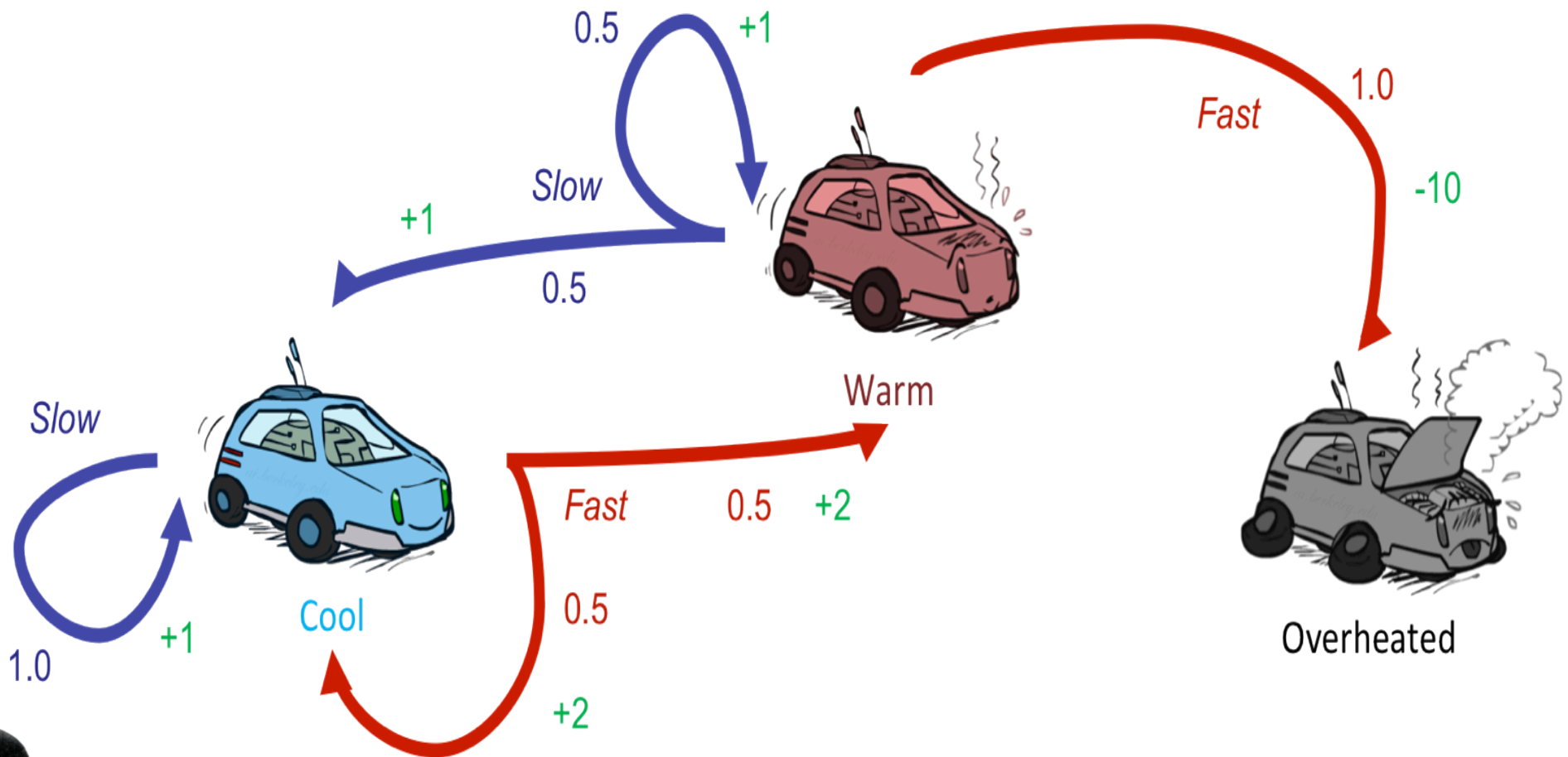
3. Extract Policy

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



Calculate $V_2(\text{warmCar})$

Assume $\gamma=1$



Slide adapted from Klein and Abbeel

Carnegie Mellon University

For General Practice, check can calculate $V_2(\text{warmCar})$



V_2

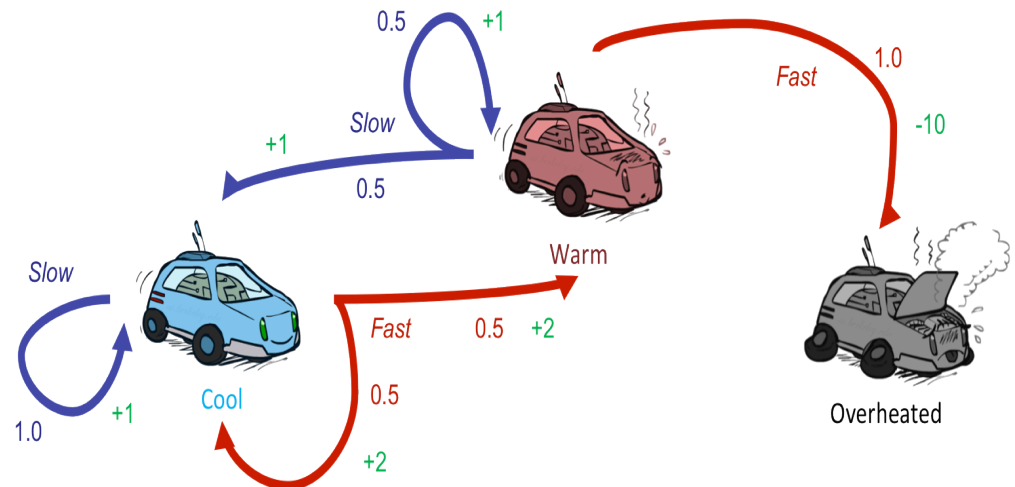
--	--	--	--

V_1

2	1	0
---	---	---

V_0

0	0	0
---	---	---



Assume $\gamma=1$

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$



Value Iteration

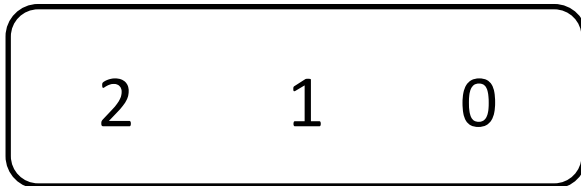
V_2 (): 2.5



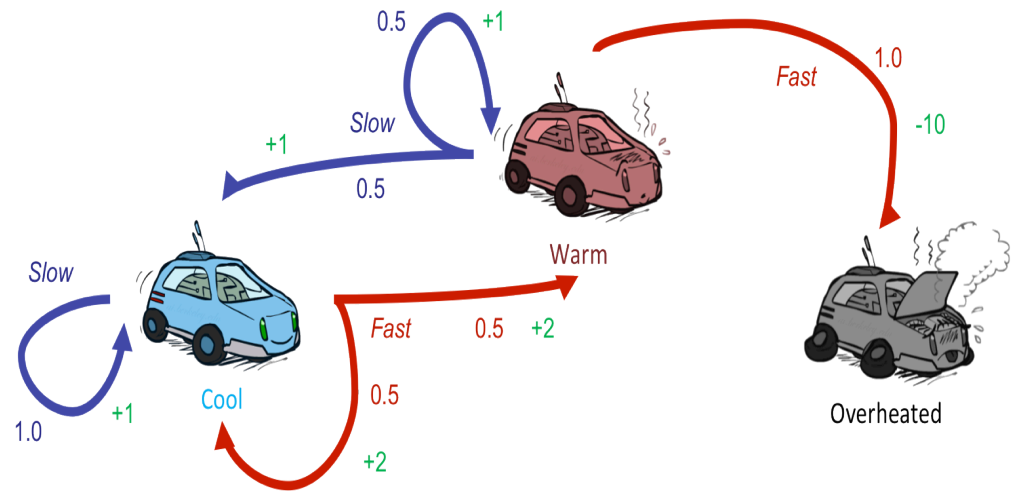
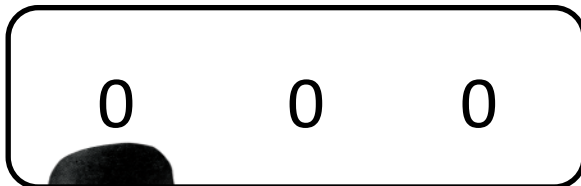
V_2



V_1



V_0



Assume $\gamma=1$

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$



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

- For all s ,

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$



COMPUTATIONAL COST PER ITERATION?

AS²

- For all s ,

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$



WILL VALUE ITERATION CONVERGE FOR INFINITE HORIZON PROBLEMS?



CONTRACTION OPERATOR

- Let O be an operator
- If $|OV - OV'| \leq |V - V'|$
- Then O is a contraction operator



WILL VALUE ITERATION CONVERGE?

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



BELLMAN OPERATOR IS A CONTRACTION ($\gamma < 1$)

$\|V - V'\|$ = Infinity norm (find max difference over all states, e.g. $\max(s) |V(s) - V'(s)|$)

$$\begin{aligned}\|BV - BV'\| &= \left\| \max_a \left[R(s, a) + \gamma \sum_{s_j \in S} p(s_j | s_i, a) V(s_j) \right] - \max_{a'} \left[R(s, a') + \gamma \sum_{s_j \in S} p(s_j | s_i, a') V'(s_j) \right] \right\| \\ &\leq \max_a \left\| \left[R(s, a) + \gamma \sum_{s_j \in S} p(s_j | s_i, a) V(s_j) - R(s, a) + \gamma \sum_{s_j \in S} p(s_j | s_i, a) V'(s_j) \right] \right\| \\ &\leq \gamma \max_a \left\| \left[\sum_{s_j \in S} p(s_j | s_i, a) V(s_j) - \sum_{s_j \in S} p(s_j | s_i, a) V'(s_j) \right] \right\| \\ &= \gamma \max_a \left\| \left[\sum_{s_j \in S} p(s_j | s_i, a) (V(s_j) - V'(s_j)) \right] \right\| \\ &\leq \gamma \max_{a, s_i} \sum_{s_j \in S} p(s_j | s_i, a) |V(s_j) - V'(s_j)| \\ &\leq \gamma \max_{a, s_i} \sum_{s_j \in S} p(s_j | s_i, a) \|V - V'\| \\ &= \gamma \|V - V'\|\end{aligned}$$



PROPERTIES OF CONTRACTION

- Only has 1 fixed point (the point reached 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

- Proof: $\lim_{k \rightarrow \infty} V_k = V^*$

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



DISCUSS AND REPORT BACK: DOES INITIALIZATION IMPACT FINAL VALUE?

Value Iteration Algorithm

1. 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
 - o For all s ,

$$V_k(s_i) = \max_a \left(\sum_{s_j \in S} p(s_j | s_i, a) [R(s, a, s') + \gamma V_{k-1}(s_j)] \right)$$

