# Reinforcement Learning



State: s
Reward: r

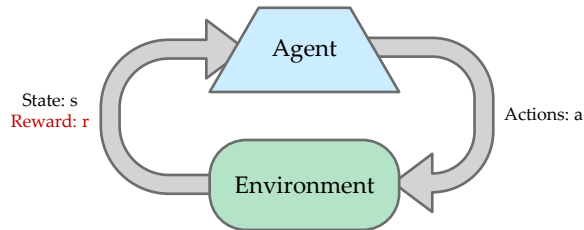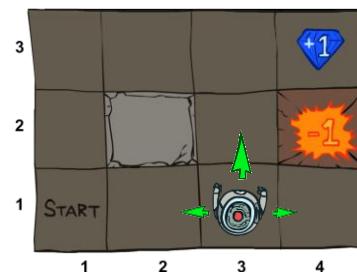Actions: a

- Basic idea:
  - Receive feedback in the form of rewards
  - Agent's utility is defined by the reward function
  - Must (learn to) act to maximize expected rewards
  - All learning is based on observed samples of outcomes!

CAL POLY
SAN LUIS OBISPO

Computer Science Department

1

1

# Example: Grid World

- A maze-like problem
  - The agent lives in a grid
  - Walls block the agent's path
- Noisy movement: actions do not always go as planned
  - 80% of the time, the 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 been taken, the agent stays put
- The agent receives rewards
  - Small "living" reward each step (can be negative)
  - Big rewards come at the end (good or bad)
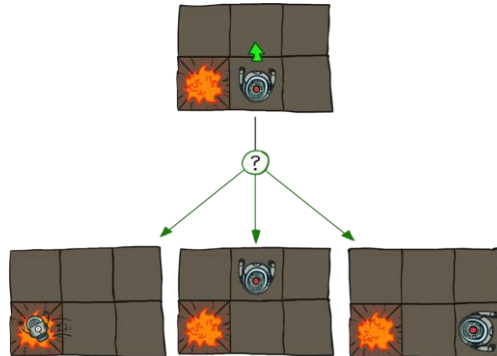- Goal: maximize sum of rewards



CAL POLY
SAN LUIS OBISPO

Computer Science Department

2

2

# Grid World Actions

Deterministic Grid World          Stochastic Grid World
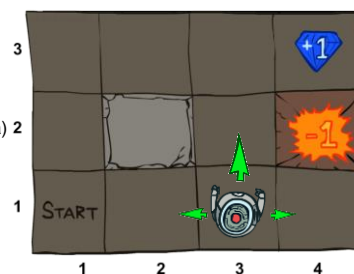
Computer Science Department                    3

3

# Markov Decision Processes

- An MDP is defined by:
    - A set of states s ∈ S
    - A set of actions a ∈ A
    - A transition function T(s, a, s')
        » Probability that a from s leads to s', i.e., Pr(s'| s, a)
        » Also called the model or the dynamics
    - A reward function R(s, a, s')
        » Sometimes just R(s) or R(s')
    - A start state
    - Possibly one or more terminal states
    - Possibly a discount factor $\gamma$ (gamma)

- MDP's are non-deterministic search problems
    - One way to solve is with expectimax search

[Demo – gridworld manual intro (L8D1)]

Computer Science Department                    4

4

## What is Markov about MDPs?

- "Markov" generally means that given the present state, the future and the past are independent

- For Markov decision processes, "Markov" means action outcomes depend only on the current state

$$P(S_{t+1} = s'|S_t = s_t, A_t = a_t, S_{t-1} = s_{t-1}, A_{t-1}, \ldots S_0 = s_0)$$
$$=$$
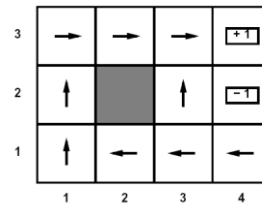$$P(S_{t+1} = s'|S_t = s_t, A_t = a_t)$$

Andrey Markov
(1856-1922)

- This is just like search, where the successor function could only depend on the current state (not the history)

**CAL POLY**
SAN LUIS OBISPO
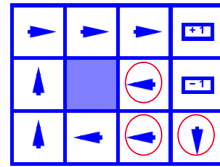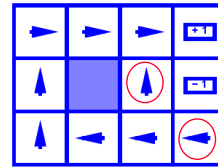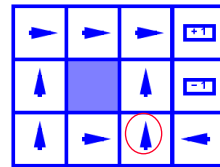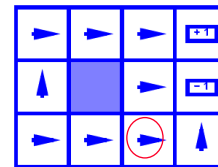
Computer Science Department

5

5

## Policies

- In deterministic single-agent search problems, we wanted an optimal plan, or sequence of actions, from start to a goal

- For MDPs, we want an optimal
  policy $\pi^*$: S → A
  - A policy $\pi$ gives an action for each state
  - An optimal policy is one that maximizes expected utility if followed
  - An explicit policy defines a reflex agent

Optimal policy when
R(s, a, s') = -0.03 for all
non-terminals s

**CAL POLY**
SAN LUIS OBISPO

Computer Science Department

6

6

# Optimal Policies for different R's



R(s) = -0.01

R(s) = -0.03

R(s) = -0.4

R(s) = -2.0

What if R(s) is > 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

7

7

# Example: Racing

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



CAL POLY
SAN LUIS OBISPO

Computer Science Department

8

8

## Racing Search Tree
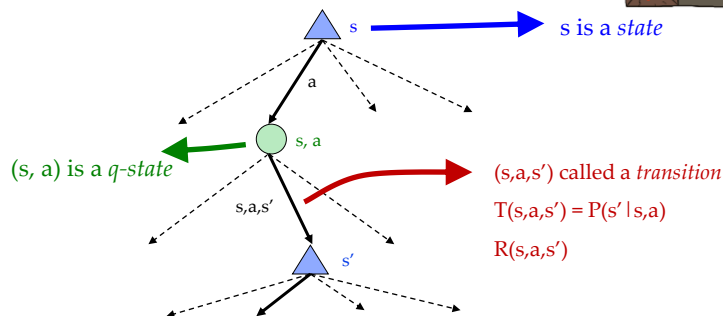
## MDP Search Trees

- Each MDP state projects an expectimax-like search tree



s is a *state*

(s, a) is a *q-state*

(s,a,s') called a *transition*

$T(s,a,s') = P(s'|s,a)$

$R(s,a,s')$

## 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]

## Discounting

- It's reasonable to maximize the sum of rewards
- It's also reasonable to prefer rewards now to rewards later
- One solution: values of rewards decay exponentially

$1$                     $\gamma$                     $\gamma^2$
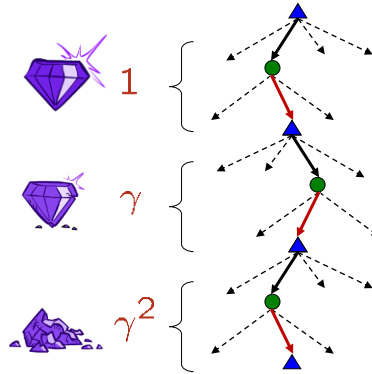
Worth Now          Worth Next Step          Worth In Two Steps

## Discounting

- How to discount?
  - Each time we descend a level, we multiply in the discount once

- Why discount?
  - Reward now is better than later
  - Can also think of it as a 1-gamma chance of ending the process at every step
  - Also helps our algorithms converge

- Example: discount of 0.5
  - U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
  - U([1,2,3]) < U([3,2,1])

$1$

$\gamma$

$\gamma^2$

CAL POLY
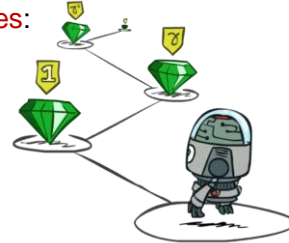SAN LUIS OBISPO

Computer Science Department

13

## Discounting

- How to discount?
  - Each time we descend a level, we multiply in the discount once

- Why discount?
  - Reward now is better than later
  - Can also think of it as a 1-gamma chance of ending the process at every step
  - Also helps our algorithms converge

- Example: discount of 0.5
  - U([1,2,3]) = 1*1 + 0.5*2 + 0.25*3
  - U([1,2,3]) < U([3,2,1])

CAL POLY
SAN LUIS OBISPO

Computer Science Department

14

14

## Stationary Preferences  -  *rational behavior Not always?*

- Theorem: if we assume stationary preferences:

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

- Then: there are only two ways to define utilities
  - Additive utility: $U([r_0, r_1, r_2, \ldots]) = r_0 + r_1 + r_2 + \cdots$
  - Discounted utilit $U([r_0, r_1, r_2, \ldots]) = r_0 + \gamma r_1 + \gamma^2 r_2 \cdots$

CAL POLY
SAN LUIS OBISPO
Computer Science Department
15

15

## Quiz: Discounting

- Given:

| 10 | | | | 1 |
|----|--|--|--|---|
| a | b | c | d | e |

  - Actions: East, West, and Exit (only available in exit states a, e)
  - Transitions: deterministic

| 10 | ← | ← | ← | 1 |

- Quiz 1: For $\gamma = 1$, what is the optimal policy?
- Quiz 2: For $\gamma = 0.1$, what is the optimal policy?

| 10 | ← | ← | → | 1 |

- Quiz 3: For which $\gamma$ are West and East equally good when in state d?

$1\gamma = 10 \gamma^3$

CAL POLY
SAN LUIS OBISPO
Computer Science Department
16

16

8

## 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. life)
    - » Gives nonstationary policies ($\pi$ depends on time left)
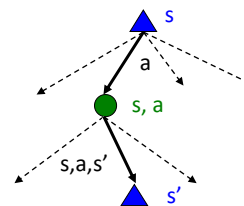  - Discounting: use $0 < \gamma < 1$

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

    - » Smaller $\gamma$ means smaller "horizon" – shorter term focus

  - Absorbing state: guarantee that for every policy, a terminal state will eventually be reached (like "overheated" for racing)

CAL POLY
SAN LUIS OBISPO

Computer Science Department

17

17

## 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'))
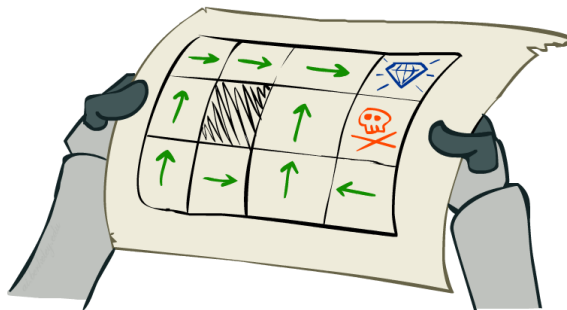  - Rewards R(s,a,s') (and discount $\gamma$)

- MDP quantities so far:
  - Policy = Choice of action for each state
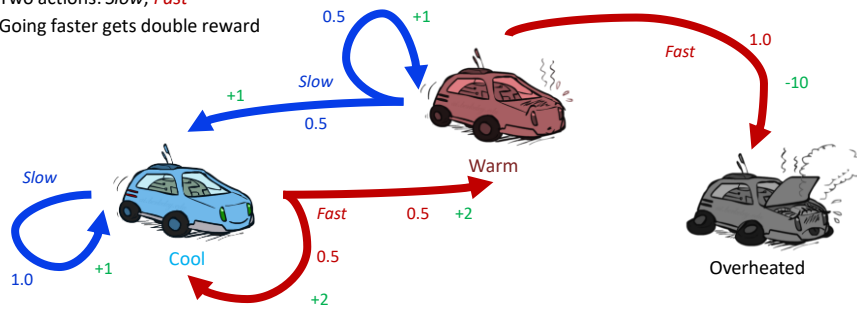  - Utility = sum of (discounted) rewards

s
a
s, a
s,a,s'
s'

CAL POLY
SAN LUIS OBISPO

Computer Science Department

18

18

# Break?  How to solve

Computer Science Department

19

19

# Solving MDPs

Computer Science Department

20

20

5/15/2024
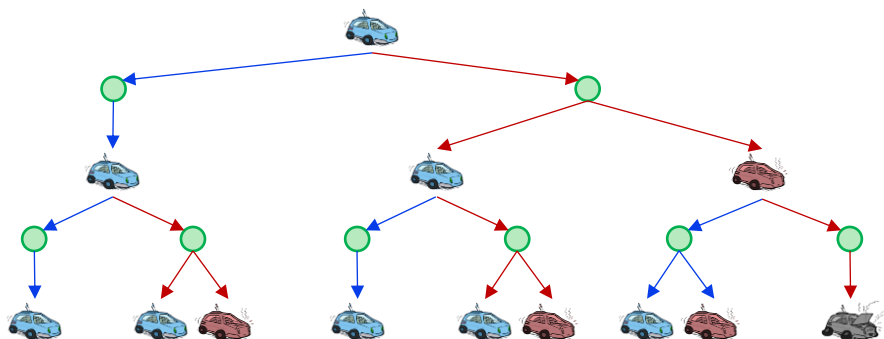
## Recall: Racing MDP

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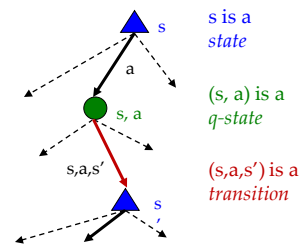- Two actions: *Slow*, *Fast*
- Going faster gets double reward



CAL POLY
SAN LUIS OBISPO

Computer Science Department

21

21

## Racing Search Tree



CAL POLY
SAN LUIS OBISPO

Computer Science Department

22

22

11

## Optimal Quantities

- The value (utility) of a state s:
  **V*(s) = expected utility** starting in s and acting optimally

- The value (utility) of a q-state (s,a):
  **Q*(s,a) = expected utility** starting out having taken action a from state s and (thereafter) acting optimally

- The optimal policy:
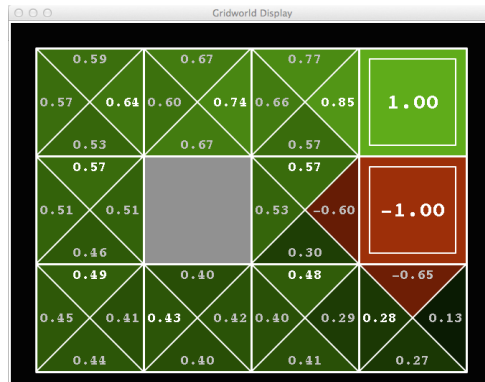  $\pi^*$(s) = optimal **action** from state s

s is a *state*

(s, a) is a *q-state*

(s,a,s') is a *transition*

Computer Science Department

23

23

## Gridworld V* Values



| | | | |
|---|---|---|---|
| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| 0.57 | | 0.57 | −1.00 |
| 0.49 | ◀ 0.43 | 0.48 | ◀ 0.28 |

Noise = 0.2
Discount = 0.9
Living reward = 0

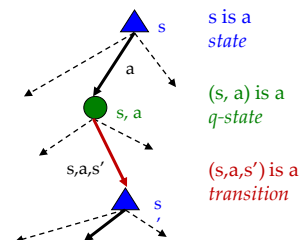Computer Science Department

24

24

12

# Gridworld Q* Values



Noise = 0.2
Discount = 0.9
Living reward = 0

# Optimal Quantities

- The value (utility) of a state s:
  **V*(s) = expected utility** starting in s and acting optimally

- The value (utility) of a q-state (s,a):
  **Q*(s,a) = expected utility** starting out having taken action a from state s and (thereafter) acting optimally

- The optimal policy:
  $\pi^*$(s) = optimal **action** from state s

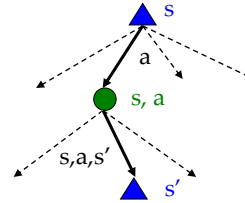s is a *state*

(s, a) is a *q-state*

(s,a,s') is a *transition*

## Values of States: Bellman Equations (Dynamic Programming)

Recursive definition of optimal value:

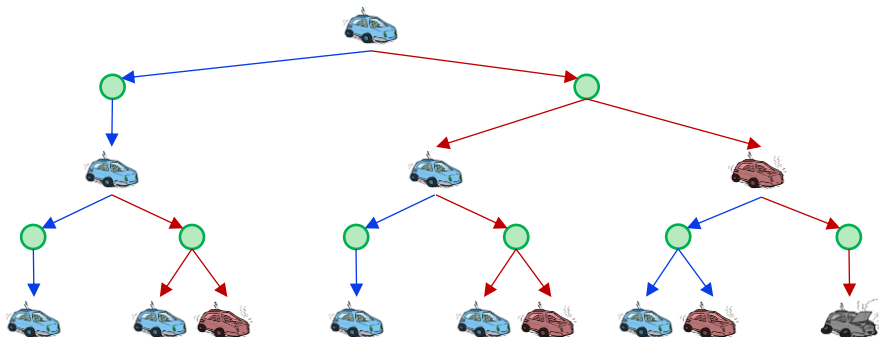$$V^*(s) = \max_a \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma V^*(s')]$$



Bellman Equations: 1 step ahead equations define optimality

$$V^*(s) = \max_a Q^*(s,a)$$
$$Q^*(s,a) = \sum_{s'} T(s,a,s')[R(s,a,s') + \gamma V^*(s')]$$

CAL POLY
SAN LUIS OBISPO                Computer Science Department                27

27

## Racing Search Tree



CAL POLY
SAN LUIS OBISPO                Computer Science Department                28

28

# Racing Search Tree

Computer Science Department

29

29

# Racing Search Tree

- We're doing way too much work with expectimax!

- Problem: States are repeated
  - Idea: Only compute needed quantities once

- Problem: Tree goes on forever
  - Idea: Do a depth-limited computation, but with increasing depths until change is small
  - Note: deep parts of the tree eventually don't matter if $\gamma < 1$

Computer Science Department

30

30

15
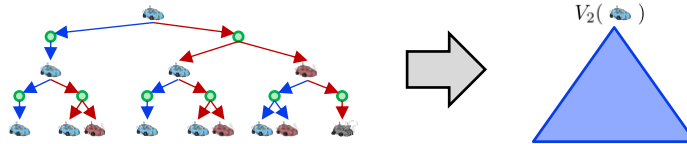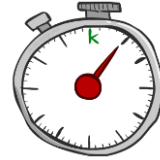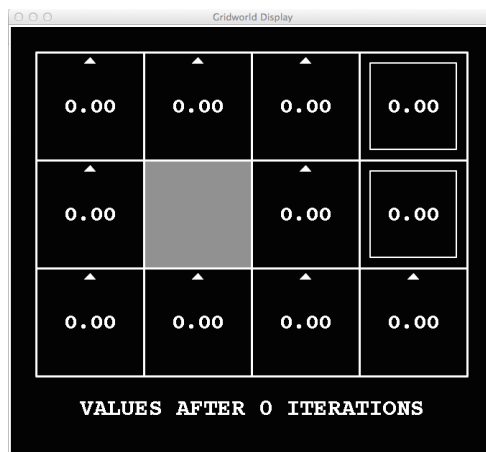
## Time-Limited Values

- Key idea: time-limited values

- Define $V_k(s)$ to be the optimal value of s if the game ends in k more time steps
  - Equivalently, it's what a depth-k expectimax would give from s

$V_2( \; )$

CAL POLY
SAN LUIS OBISPO

Computer Science Department

31

31

## k=0

Gridworld Display

| 0.00 | 0.00 | 0.00 | 0.00 |
| 0.00 |  | 0.00 | 0.00 |
| 0.00 | 0.00 | 0.00 | 0.00 |

**VALUES AFTER 0 ITERATIONS**

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

32

32

## k=1



VALUES AFTER 1 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

33

33

## k=2



VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

34

34

# k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

35

# k=4



VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

36

# k=5



VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

37

37

# k=6



VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

38

38

# k=7



VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

39

39

# k=8



VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

40

40

# k=9



VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

41

41

# k=10



VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

42

42

## k=11



VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

43

43

## k=12



VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

44

44

# k=100



Noise = 0.2
Discount = 0.9
Living reward = 0

CAL POLY
SAN LUIS OBISPO

Computer Science Department

45

# Computing Time-Limited Values



$V_4(\;)\quad V_4(\;)\quad V_4(\;)$

$V_3(\;)\quad V_3(\;)\quad V_3(\;)$

$V_2(\;)\quad V_2(\;)\quad V_2(\;)$

$V_1(\;)\quad V_1(\;)\quad V_1(\;)$

$V_0(\;)\quad V_0(\;)\quad V_0(\;)$

CAL POLY
SAN LUIS OBISPO

Computer Science Department

46

# Value Iteration

# Value Iteration

- Start with $V_0(s) = 0$: no time steps left means an expected reward sum of zero

- Given vector of $V_k(s)$ values, do one ply of expectimax from each state:

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma\, V_k(s') \right]$$



- Repeat until convergence, which yields V*

- Complexity of each iteration: $O(S^2 A)$

- Theorem: will converge to unique optimal values
  – Basic idea: approximations get refined towards optimal values
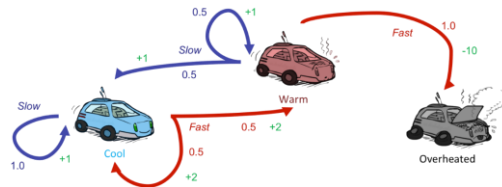  – Policy may converge long before values do

# Example: Value Iteration



$V_2$

$V_1$ | S: 1
F: .5*2+.5*2=2

$V_0$ | 0      0      0

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$

# Example: Value Iteration



$V_2$

$V_1$ | 2 | S: .5*1+.5*1=1
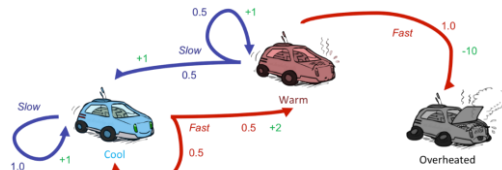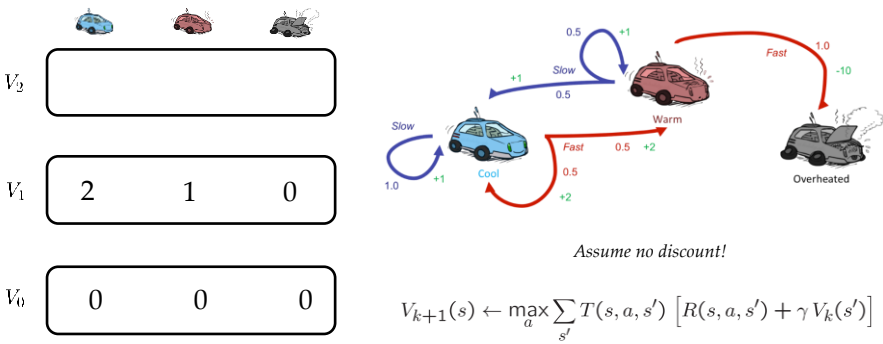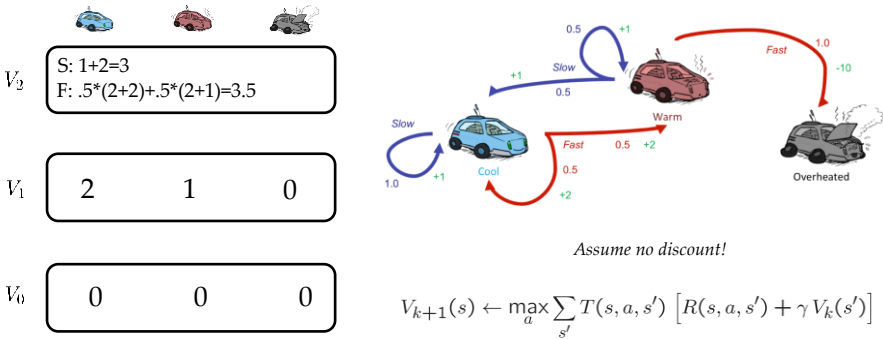F: -10

$V_0$ | 0      0      0

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$

## Example: Value Iteration



$V_2$

$V_1$ | 2 | 1 | 0

$V_0$ | 0 | 0 | 0

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$

CAL POLY
SAN LUIS OBISPO

Computer Science Department

51

51

## Example: Value Iteration



$V_2$ | S: 1+2=3
F: .5*(2+2)+.5*(2+1)=3.5

$V_1$ | 2 | 1 | 0

$V_0$ | 0 | 0 | 0

*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$

CAL POLY
SAN LUIS OBISPO

Computer Science Department

52

52

# Example: Value Iteration



*Assume no discount!*

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s,a,s') \left[ R(s,a,s') + \gamma V_k(s') \right]$$
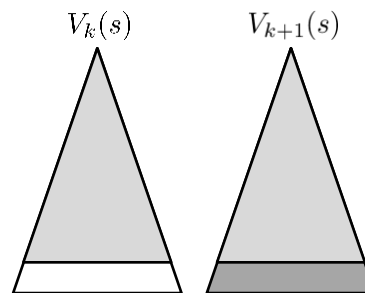
| $V_2$ | 3.5 | 2.5 | 0 |
| $V_1$ | 2 | 1 | 0 |
| $V_0$ | 0 | 0 | 0 |

53

# Convergence*

- How do we know the $V_k$ vectors are going to converge? (assuming $0 < \gamma < 1$)

- Proof Sketch:
    - For any state $V_k$ and $V_{k+1}$ can be viewed as depth k+1 expectimax results in nearly identical search trees
    - The difference is that on the bottom layer, $V_{k+1}$ has actual rewards while $V_k$ has zeros
    - That last layer is at best all $R_{MAX}$
    - It is at worst $R_{MIN}$
    - But everything is discounted by $\gamma^k$ that far out
    - So $V_k$ and $V_{k+1}$ are at most $\gamma^k \max|R|$ different
    - So as k increases, the values converge

$V_k(s)$    $V_{k+1}(s)$

54

27

# Next Lecture: Policy-Based Methods

CAL POLY
SAN LUIS OBISPO

Computer Science Department

55

55