# CS 573: Artificial Intelligence

## Markov Decision Processes
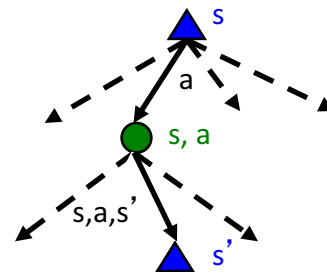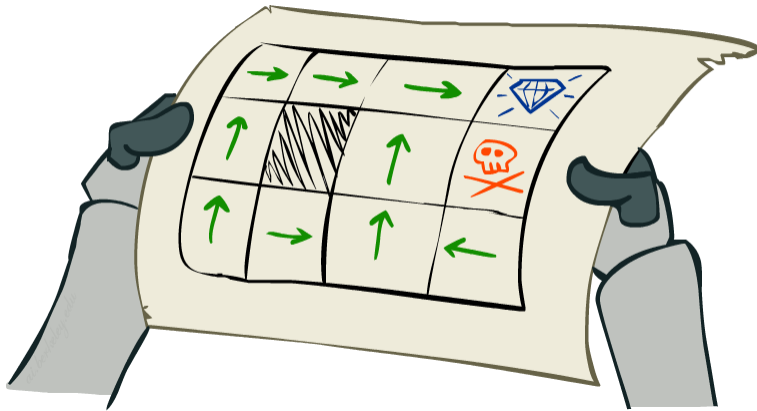


Dan Weld

University of Washington

---

# 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

# Solving MDPs



- Value Iteration
  - Asynchronous VI

- Policy Iteration

- Reinforcement Learning

# V* = Optimal Value Function

The value (utility) of a state s:

$$V^*(s)$$

"expected utility starting in s & acting optimally forever"

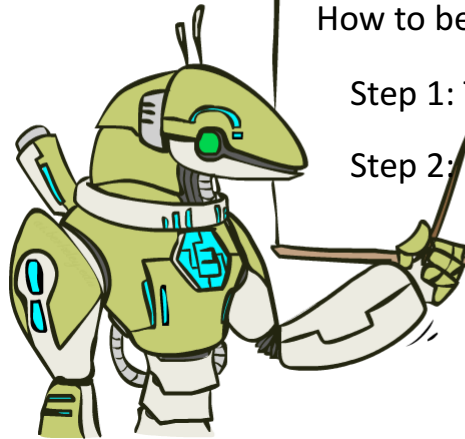## Q*

The value (utility) of the q-state (s,a):

$$Q^*(s,a)$$

"expected utility of 1) starting in state **s**
2) taking action **a**
3) acting *optimally* forever after that"

Q*(s,a) = reward from executing a in s then ending in s'
plus… discounted value of V*(s')

## $\pi^*$    Specifies The Optimal Policy

$\pi^*(s)$ = optimal action from state s

# The Bellman Equations



How to be optimal:

Step 1: Take correct first action

Step 2: Keep being optimal

# The Bellman Equations



- Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values
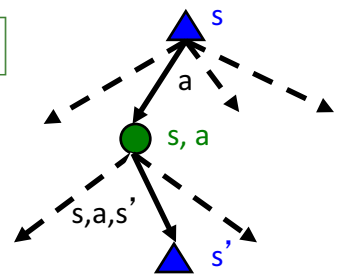
(1920-1984)

$$V^*(s) = \max_a Q^*(s, a)$$

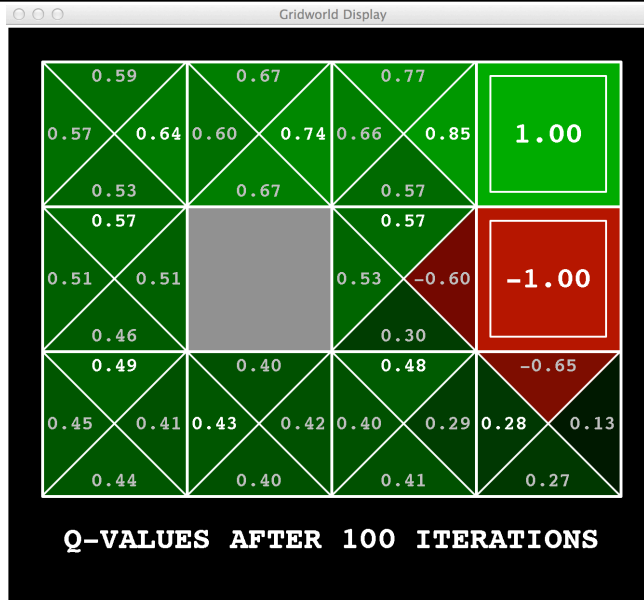$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

- These are the Bellman equations, and they characterize optimal values in a way we'll use over and over

# Gridworld: Q*



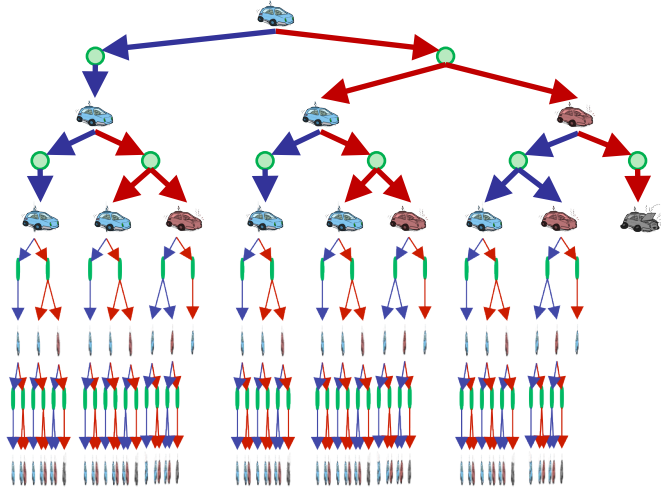Q-VALUES AFTER 100 ITERATIONS

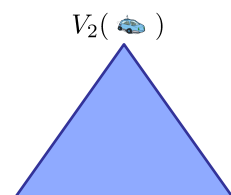# Gridworld Values V* $\quad V^*(s) = \max_a Q^*(s, a)$
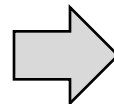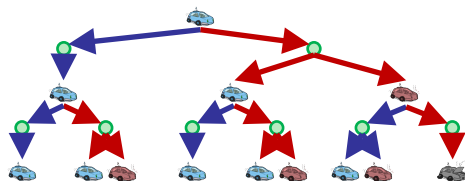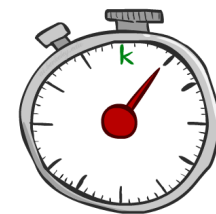


VALUES AFTER 100 ITERATIONS

# No End in Sight…

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

- Problem 1: States are repeated
  - Idea: Only compute needed quantities once
  - Like **graph search** (*vs.* tree search)

- Problem 2: Tree goes on forever
  - Rewards @ each step → **V changes**
  - Idea: Do a depth-limited computation, but with increasing depths until change is small
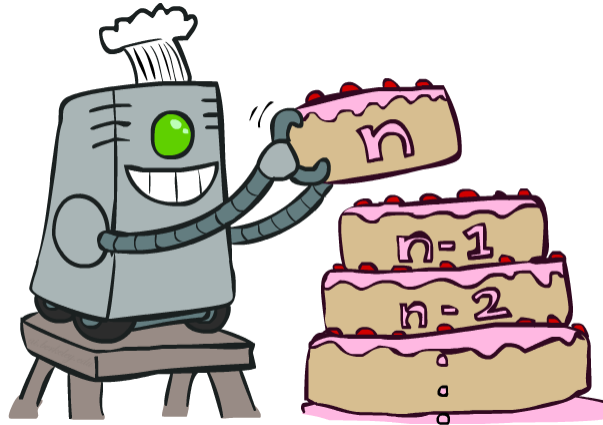  - Note: deep parts of the tree *eventually* don't matter if $\gamma < 1$



# 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



$V_2( \text{🚗} )$

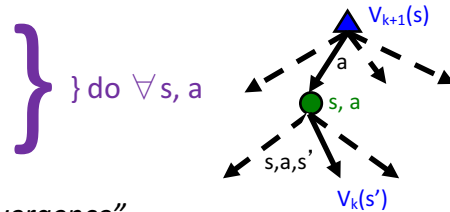# Value Iteration



---

# Value Iteration

- **Forall s, initialize $V_0(s) = 0$**    *no time steps left means an expected reward of zero*
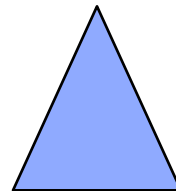
- **Repeat**
    K += 1

    $Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$    } do $\forall$ s, a

    $V_{k+1}(s) = Max_a Q_{k+1}(s, a)$

- **Repeat until $|V_{k+1}(s) - V_k(s)| < \varepsilon$,    forall s**    *"convergence"*

    Successive approximation; dynamic programming

$V_{k+1}(s)$
a
s, a
s,a,s'
$V_k(s')$

# Example: Value Iteration

*Assume no discount (gamma=1) to keep math simple!*



$$Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$$

$$V_{k+1}(s) = \text{Max}_a Q_{k+1}(s, a)$$

---

# Example: Value Iteration
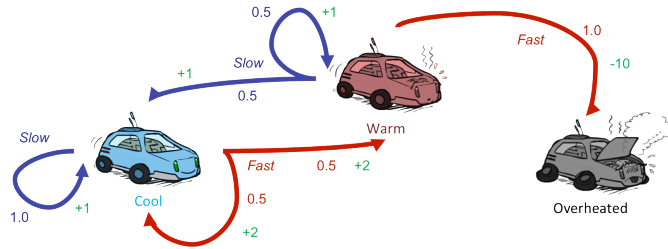
*Assume no discount (gamma=1) to keep math simple!*

$V_0$ | 0 | 0 | 0

$V_1$

$V_2$
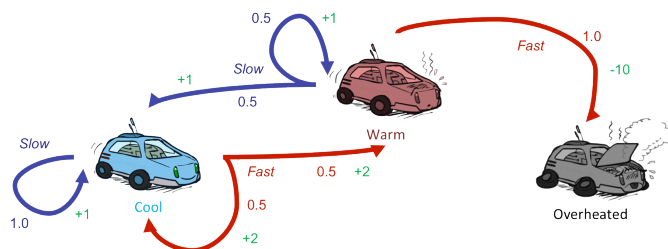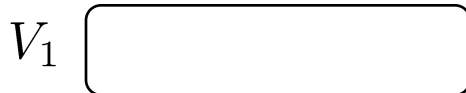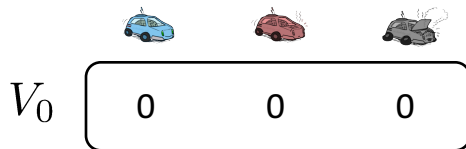


$$Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$$

$$V_{k+1}(s) = \text{Max}_a Q_{k+1}(s, a)$$

Example: Value Iteration

Q( 🚗 ,fast) =

Q( 🚗 slow) =

math simple!

$V_0$ | 0 | 0 | 0

$Q_1(s,a)=$ 0

$V_1$ | | 0

$Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$

$V_{k+1}(s) = Max_a Q_{k+1}(s, a)$

$V_2$ |

---



Example: Value Iteration

Q( 🚗 ,fast) = -10 + 0

Q( 🚗 slow) = ½(1 + 0) + ½(1+0)

math simple!

$V_0$ | 0 | 0 | 0

$Q_1(s,a)=$ 1, -10 0

$V_1$ | 1 | 0

$Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$

$V_{k+1}(s) = Max_a Q_{k+1}(s, a)$

$V_2$ |

# Example: Value Iteration

*...ma=1) to keep math simple!*

$V_0$

| 0 | 0 | 0 |

$Q_1(s,a)=$   1, 2     1,-10     0

$V_1$

| 2 | 1 | 0 |

$Q_{k+1}(s, a) = \Sigma_{s'}\ T(s, a, s') [ R(s, a, s') + \gamma\ V_k(s')]$
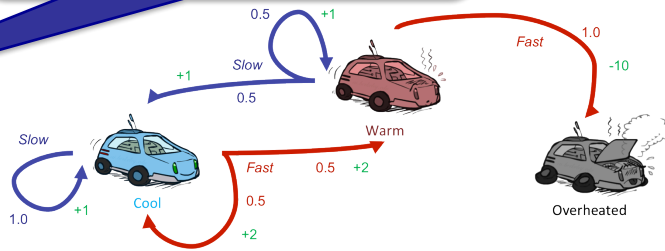
$V_{k+1}(s) = \text{Max}_a\ Q_{k+1}(s, a)$

$V_2$

| | | |

---

# Example: Value Iteration

*Assume no discount (gamma=1) to keep math simple!*

$V_0$

| 0 | 0 | 0 |

$Q_1(s,a)=$   1, 2     1,-10     0

$V_1$

| 2 | 1 | 0 |

$Q_2(s,a)=$   3,3.5     2.5,-10     0

$V_2$

| 3.5 | 2.5 | 0 |

$Q_{k+1}(s, a) = \Sigma_{s'}\ T(s, a, s') [ R(s, a, s') + \gamma\ V_k(s')]$
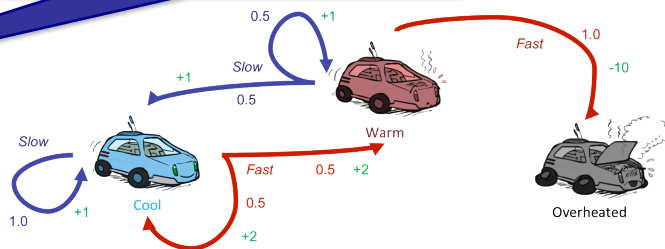
$V_{k+1}(s) = \text{Max}_a\ Q_{k+1}(s, a)$

# k=0



**VALUES AFTER 1 ITERATIONS**

Noise = 0.2
Discount = 0.9
Living reward = 0

---

# k=1

If agent is in 4,3, it only has one legal action: get jewel. It gets a reward and the game is over.
If agent is in the pit, it has only one legal action, die. It gets a penalty and the game is over.

Agent does NOT get a reward for moving INTO 4,3.



**VALUES AFTER 1 ITERATIONS**

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=2



VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=4



VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=5



VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=6



VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=7



VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=8



VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=9



VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=10



VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=11



VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=12



VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# k=100



VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

# VI: Policy Extraction



# Computing Actions from Values

- Let's imagine we have the optimal values V*(s)

- How should we act?
  - In general, it's not obvious!

- We need to do a mini-expectimax (one step)

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

- This is called policy extraction, since it gets the policy implied by the values

# Computing Actions from Q-Values

- Let's imagine we have the optimal q-values:

- How should we act?
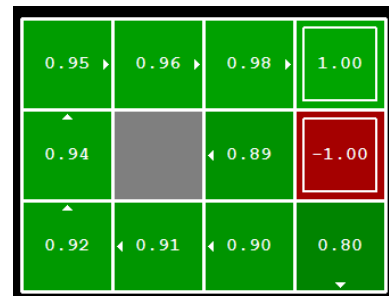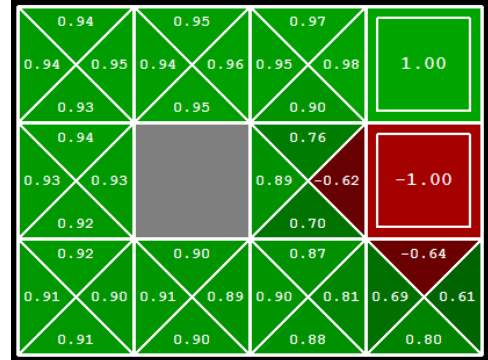  - Completely trivial to decide!

$$\pi^*(s) = \arg\max_a Q^*(s, a)$$



| | | | |
|---|---|---|---|
| 0.94 | 0.95 | 0.97 | |
| 0.94 / 0.95 | 0.94 / 0.96 | 0.95 / 0.98 | 1.00 |
| 0.93 | 0.95 | 0.90 | |
| 0.94 | | 0.76 | |
| 0.93 / 0.93 | | 0.89 / -0.62 | -1.00 |
| 0.92 | | 0.70 | |
| 0.92 | 0.90 | 0.87 | -0.64 |
| 0.91 / 0.90 | 0.91 / 0.89 | 0.90 / 0.81 | 0.69 / 0.61 |
| 0.91 | 0.90 | 0.88 | 0.80 |

- Important lesson: actions are easier to select from q-values than values!

---

# Value Iteration - Recap

- **Forall s, Initialize $V_0(s) = 0$**    *no time steps left means an expected reward of zero*

- **Repeat**                 *do Bellman backups*
  K += 1
  Repeat for all states, s, and all actions, a:

  $Q_{k+1}(s, a) = \Sigma_{s'} T(s, a, s') [ R(s, a, s') + \gamma V_k(s')]$

  $V_{k+1}(s) = \text{Max}_a Q_{k+1}(s, a)$

- **Until $|V_{k+1}(s) - V_k(s)| < \varepsilon$,    forall s**   *"convergence"*

- **Theorem: will converge to unique optimal values**