

CSE 573: Artificial Intelligence

Autumn 2010

Lecture 7: MDPs/RL
10/21/2010

Luke Zettlemoyer

Many slides over the course adapted from either Dan Klein,
Stuart Russell or Andrew Moore

Outline

- Markov decision processes
 - Review Optimality / Value Iteration
 - Value Iteration convergence / complexity
 - Policy Iteration
- Reinforcement Learning
 - Passive Learning
 - TD Updates
 - Q-learning
- 3:30: Tom Mitchell's Distinguished Lecture
 - EEB-105

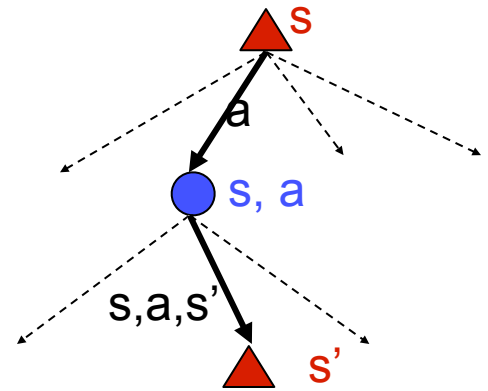
Homework Rant

- PS2 Due Tuesday!
- PS1 will be handed back this afternoon
 - Admissibility was hard, but overall everyone did well!
 - **Next time:** Follow the instructions!!!
 - Only hand in the one/two requested files (and don't zip/tar them)
 - Don't change any other files
 - Turn off your debug printouts
 - Comment your code (if you want partial credit)

Recap: MDPs

- Markov decision processes:

- States S
- Actions A
- Transitions $P(s'|s,a)$ (or $T(s,a,s')$)
- Rewards $R(s,a,s')$ (and discount γ)
- Start state s_0



- Quantities:

- Policy = map of states to actions
- Utility = sum of discounted rewards
- Values = expected future utility from a state
- Q-Values = expected future utility from a q-state

Recap: Optimal Utilities

- The utility of a state s :

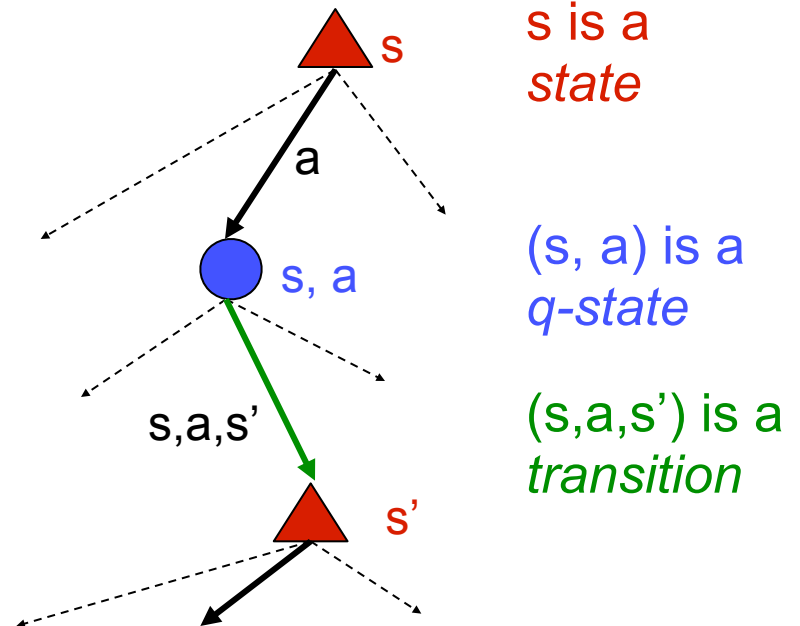
$V^*(s)$ = expected utility starting in s and acting optimally

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

$Q^*(s,a)$ = expected utility starting in s , taking action a and thereafter acting optimally

- The optimal policy:

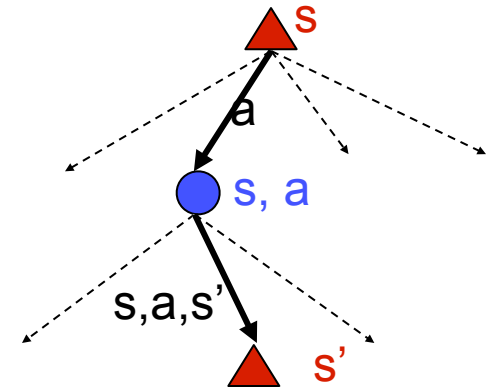
$\pi^*(s)$ = optimal action from state s



Recap: Bellman Equations

- Definition of utility leads to a simple one-step lookahead relationship amongst optimal utility values:

Total optimal rewards = maximize over choice of (first action plus optimal future)



- Formally:

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

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

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

Practice: Computing Actions

- Which action should we chose from state s :
 - Given optimal values V ?

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

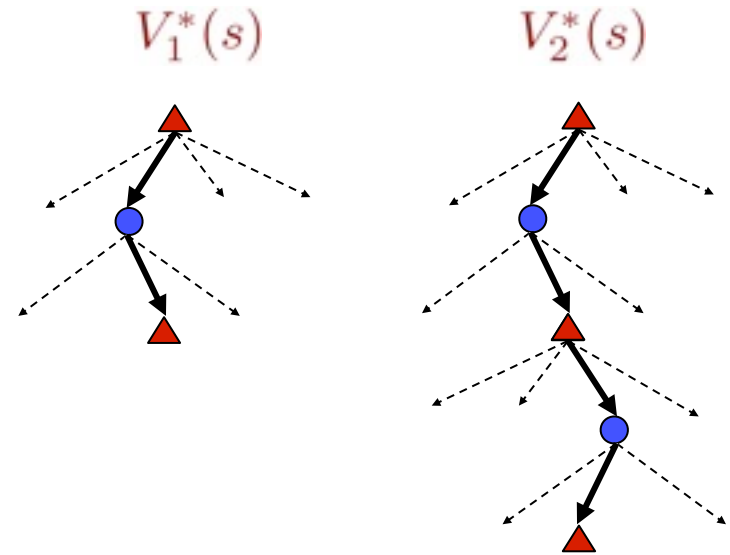
- Given optimal q-values Q ?

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

- Lesson: actions are easier to select from Q 's!

Value Estimates

- Calculate estimates $V_k^*(s)$
 - Not the optimal value of s !
 - The optimal value considering only next k time steps (k rewards)
 - As $k \rightarrow \infty$, it approaches the optimal value
- Value Iteration: dynamic programming



Example: Value Iteration

▲ 0.00	▲ 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

Value Iteration

- Idea:

- Start with $V_0^*(s) = 0$, which we know is right (why?)
- Given V_i^* , calculate the values for all states for depth $i+1$:

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

- Throw out old vector V_i^*
 - Repeat until convergence
 - This is called a **value update** or **Bellman update**
- **Theorem: will converge to unique optimal values**
 - Basic idea: approximations get refined towards optimal values
 - Policy may converge long before values do

Convergence

- Define the max-norm: $\|U\| = \max_s |U(s)|$

- Theorem: For any two value vectors U and V

$$\|U^{t+1} - V^{t+1}\| \leq \gamma \|U^t - V^t\|$$

- I.e. any distinct approximations must get closer to each other, so, in particular, any approximation must get closer to the true U and value iteration converges to a unique, stable, optimal solution

- Theorem:

$$\|U^{t+1} - U^t\| < \epsilon, \Rightarrow \|U^{t+1} - U\| < 2\epsilon\gamma/(1 - \gamma)$$

- I.e. once the change in our approximation is small, it must also be close to correct

Value Iteration Complexity

- Problem size:
 - $|A|$ actions and $|S|$ states
- Each Iteration
 - Computation: $O(|A| \cdot |S|^2)$
 - Space: $O(|S|)$
- Num of iterations
 - Can be exponential in the discount factor γ

Asynchronous Value Iteration*

- In value iteration, we update every state in each iteration
- Actually, *any* sequences of Bellman updates will converge if every state is visited infinitely often
- In fact, we can update the policy as seldom or often as we like, and we will still converge
- Idea: Update states whose value we expect to change:
If $|V_{i+1}(s) - V_i(s)|$ is large then update predecessors of s

Utilities for Fixed Policies

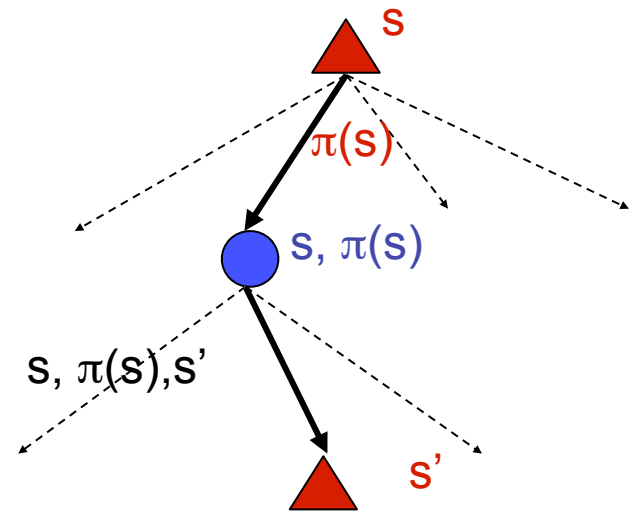
- Another basic operation:
compute the utility of a state s
under a fixed (general non-optimal)
policy

- Define the utility of a state s ,
under a fixed policy π :

$V^\pi(s)$ = expected total discounted
rewards (return) starting in s and
following π

- Recursive relation (one-step
look-ahead / Bellman equation):

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



Policy Evaluation

- How do we calculate the V 's for a fixed policy?
- Idea one: modify Bellman updates

$$V_0^\pi(s) = 0$$

$$V_{i+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_i^\pi(s')]$$

- Idea two: it's just a linear system, solve with Matlab (or whatever)

Policy Iteration

- Problem with value iteration:
 - Considering all actions each iteration is slow: takes $|A|$ times longer than policy evaluation
 - But policy doesn't change each iteration, time wasted
- Alternative to value iteration:
 - **Step 1: Policy evaluation:** calculate utilities for a fixed policy (not optimal utilities!) until convergence (fast)
 - **Step 2: Policy improvement:** update policy using one-step lookahead with resulting converged (but not optimal!) utilities (slow but infrequent)
 - Repeat steps until policy converges

Policy Iteration

- Policy evaluation: with fixed current policy π , find values with simplified Bellman updates:
 - Iterate until values converge

$$V_{i+1}^{\pi_k}(s) \leftarrow \sum_{s'} T(s, \pi_k(s), s') [R(s, \pi_k(s), s') + \gamma V_i^{\pi_k}(s')]$$

- Policy improvement: with fixed utilities, find the best action according to one-step look-ahead

$$\pi_{k+1}(s) = \arg \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^{\pi_k}(s')]$$

Policy Iteration Complexity

- Problem size:
 - $|A|$ actions and $|S|$ states
- Each Iteration
 - Computation: $O(|S|^3 + |A| \cdot |S|^2)$
 - Space: $O(|S|)$
- Num of iterations
 - Unknown, but can be faster in practice
 - Convergence is guaranteed

Comparison

- **In value iteration:**
 - Every pass (or “backup”) updates both utilities (explicitly, based on current utilities) and policy (possibly implicitly, based on current policy)
- **In policy iteration:**
 - Several passes to update utilities with frozen policy
 - Occasional passes to update policies
- **Hybrid approaches (asynchronous policy iteration):**
 - Any sequences of partial updates to either policy entries or utilities will converge if every state is visited infinitely often

What is it doing?

-

Step Delay: 0.10000

+

-

Epsilon: 0.500

+

-

Discount: 0.800

+

-

Learning Rate: 0.800

+



Step: 75

Position: 63

Velocity: -6.04

100-step Avg Velocity: 0.68

Reinforcement Learning

- Reinforcement learning:

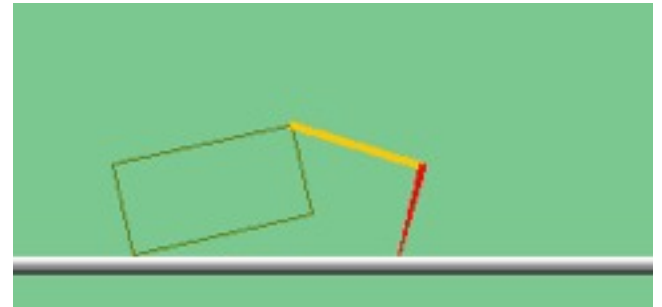
- Still have an MDP:

- A set of states $s \in S$
 - A set of actions (per state) A
 - A model $T(s,a,s')$
 - A reward function $R(s,a,s')$

- Still looking for a policy $\pi(s)$

- New twist: **don't know T or R**

- I.e. don't know which states are good or what the actions do
 - Must actually try actions and states out to learn

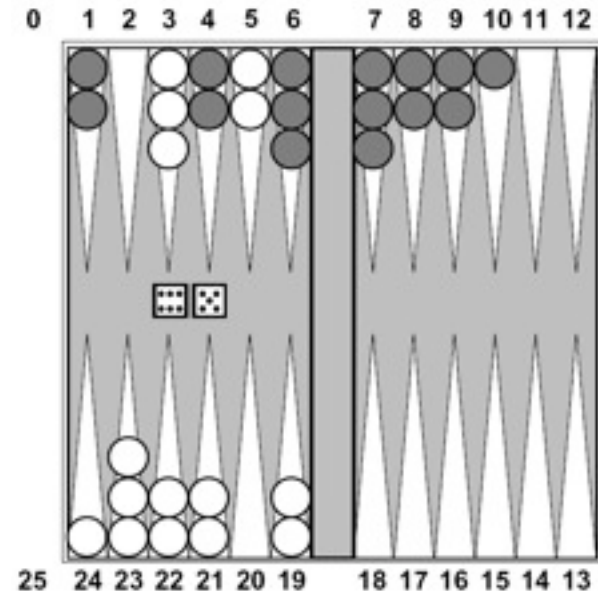


Example: Animal Learning

- RL studied experimentally for more than 60 years in psychology
 - Rewards: food, pain, hunger, drugs, etc.
 - Mechanisms and sophistication debated
- Example: foraging
 - Bees learn near-optimal foraging plan in field of artificial flowers with controlled nectar supplies
 - Bees have a direct neural connection from nectar intake measurement to motor planning area

Example: Backgammon

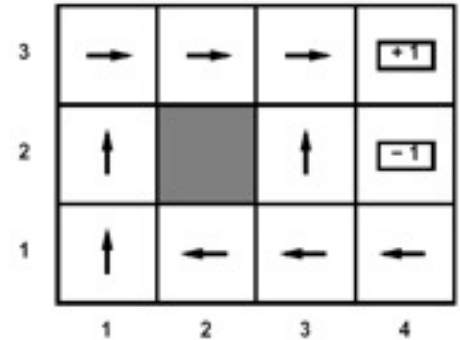
- Reward only for win / loss in terminal states, zero otherwise
- TD-Gammon learns a function approximation to $V(s)$ using a neural network
- Combined with depth 3 search, one of the top 3 players in the world
- You could imagine training Pacman this way...
- ... but it's tricky! (It's also P3)



Passive Learning

- Simplified task

- You don't know the transitions $T(s,a,s')$
- You don't know the rewards $R(s,a,s')$
- You are given a policy $\pi(s)$
- **Goal: learn the state values** (and maybe the model)
- I.e., policy evaluation



- In this case:

- Learner “along for the ride”
- No choice about what actions to take
- Just execute the policy and learn from experience
- We'll get to the active case soon
- This is NOT offline planning!

Detour: Sampling Expectations

- Want to compute an expectation weighted by $P(x)$:

$$E[f(x)] = \sum_x P(x) f(x)$$

- Model-based: estimate $P(x)$ from samples, compute expectation

$$x_i \sim P(x)$$
$$\hat{P}(x) = \text{count}(x)/k$$
$$E[f(x)] \approx \sum_x \hat{P}(x) f(x)$$

- Model-free: estimate expectation directly from samples

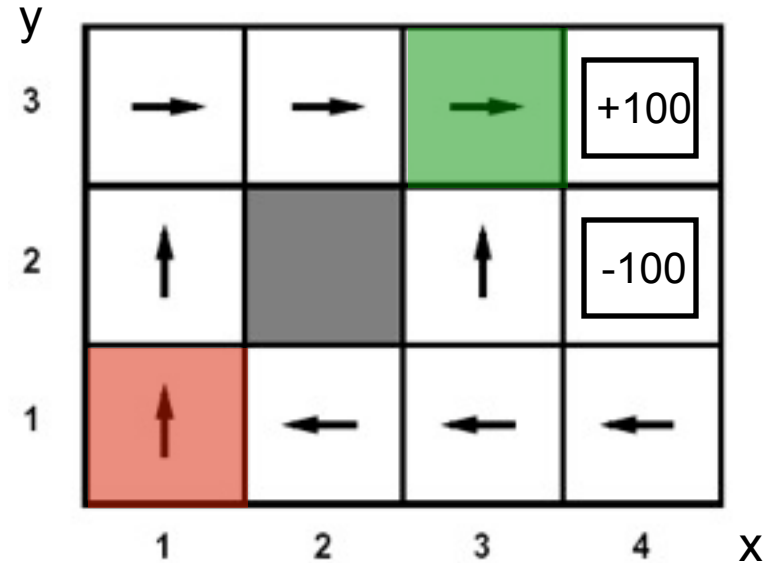
$$x_i \sim P(x)$$
$$E[f(x)] \approx \frac{1}{k} \sum_i f(x_i)$$

- Why does this work? Because samples appear with the right frequencies!

Example: Direct Estimation

Episodes:

(1,1) up -1	(1,1) up -1
(1,2) up -1	(1,2) up -1
(1,2) up -1	(1,3) right -1
(1,3) right -1	(2,3) right -1
(2,3) right -1	(3,3) right -1
(3,3) right -1	(3,2) up -1
(3,2) up -1	(4,2) exit -100
(3,3) right -1	(done)
(4,3) exit +100	
(done)	



$\gamma = 1, R = -1$

$$V(1,1) \sim (92 + -106) / 2 = -7$$

$$V(3,3) \sim (99 + 97 + -102) / 3 = 31.3$$

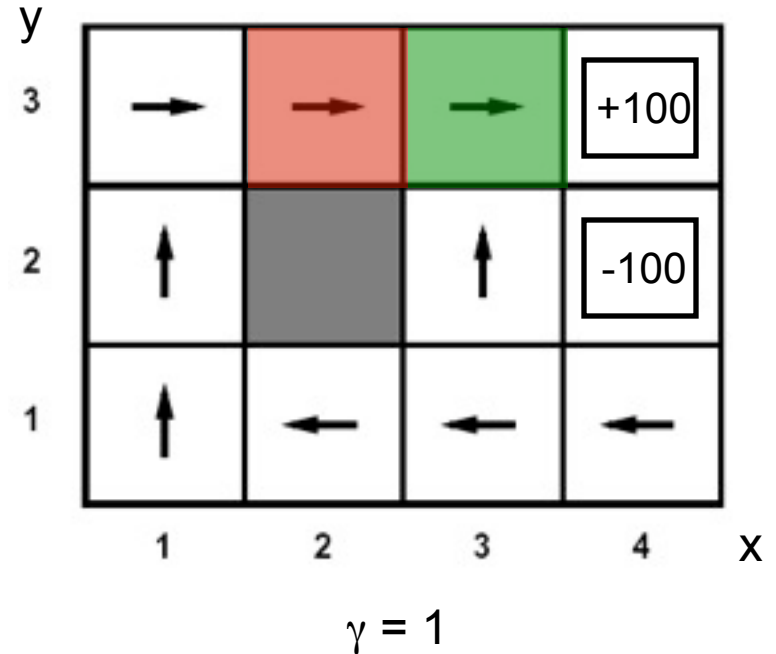
Model-Based Learning

- Idea:
 - Learn the model empirically (rather than values)
 - Solve the MDP as if the learned model were correct
 - Better than direct estimation?
- Empirical model learning
 - Simplest case:
 - Count outcomes for each s,a
 - Normalize to give estimate of $T(s,a,s')$
 - Discover $R(s,a,s')$ the first time we experience (s,a,s')
 - More complex learners are possible (e.g. if we know that all squares have related action outcomes, e.g. “stationary noise”)

Example: Model-Based Learning

Episodes:

- | | |
|-----------------|-----------------|
| (1,1) up -1 | (1,1) up -1 |
| (1,2) up -1 | (1,2) up -1 |
| (1,2) up -1 | (1,3) right -1 |
| (1,3) right -1 | (2,3) right -1 |
| (2,3) right -1 | (3,3) right -1 |
| (3,3) right -1 | (3,2) up -1 |
| (3,2) up -1 | (4,2) exit -100 |
| (3,3) right -1 | (done) |
| (4,3) exit +100 | |
| (done) | |

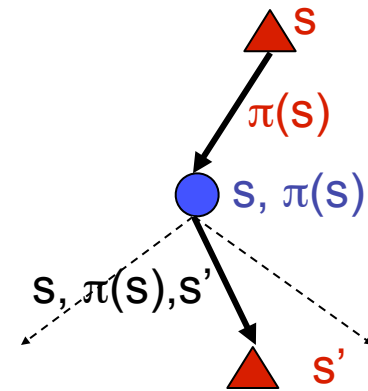


$$T(\langle 3,3 \rangle, \text{right}, \langle 4,3 \rangle) = 1 / 3$$

$$T(\langle 2,3 \rangle, \text{right}, \langle 3,3 \rangle) = 2 / 2$$

Recap: Model-Based Policy Evaluation

- Simplified Bellman updates to calculate V for a fixed policy:
 - New V is expected one-step-look-ahead using current V
 - Unfortunately, need T and R



$$V_0^\pi(s) = 0$$

$$V_{i+1}^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_i^\pi(s')]$$

Sample Avg to Replace Expectation?

$$V_{i+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_i^{\pi}(s')]$$

- Who needs T and R? Approximate the expectation with samples (drawn from T!)

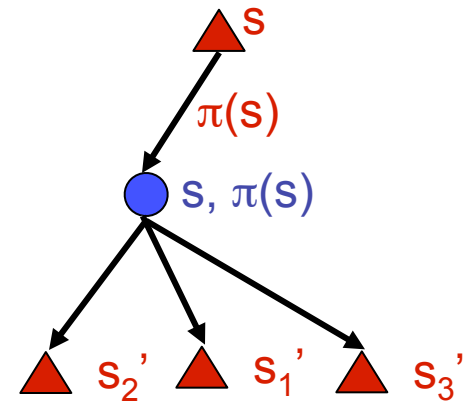
$$\text{sample}_1 = R(s, \pi(s), s'_1) + \gamma V_i^{\pi}(s'_1)$$

$$\text{sample}_2 = R(s, \pi(s), s'_2) + \gamma V_i^{\pi}(s'_2)$$

...

$$\text{sample}_k = R(s, \pi(s), s'_k) + \gamma V_i^{\pi}(s'_k)$$

$$V_{i+1}^{\pi}(s) \leftarrow \frac{1}{k} \sum_i \text{sample}_i$$



Exponential Moving Average

- Exponential moving average
 - Makes recent samples more important

$$\bar{x}_n = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- Forgets about the past (distant past values were wrong anyway)
- Easy to compute from the running average

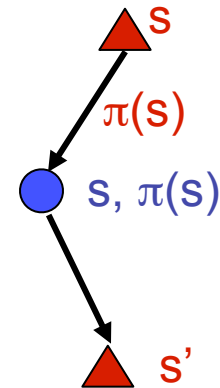
$$\bar{x}_n = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$$

- Decreasing learning rate can give converging averages

Model-Free Learning

$$V^\pi(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

- Big idea: why bother learning T?
 - Update V each time we experience a transition
- Temporal difference learning (TD)
 - Policy still fixed!
 - Move values toward value of whatever successor occurs: running average!



$$sample = R(s, \pi(s), s') + \gamma V^\pi(s')$$

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + (\alpha)sample$$

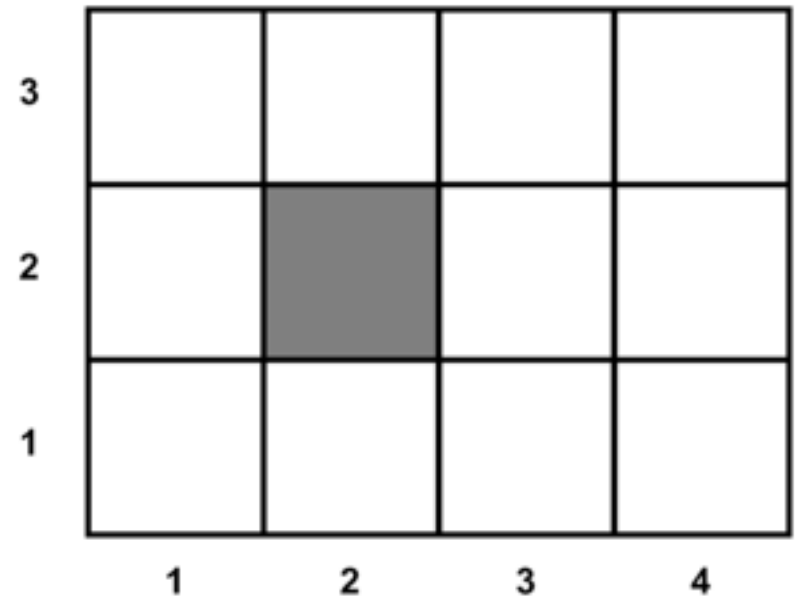
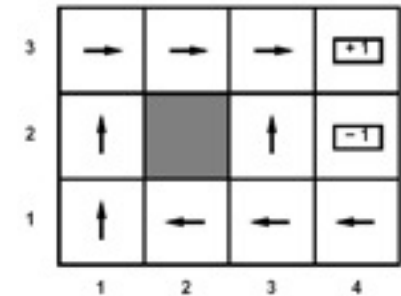
$$V^\pi(s) \leftarrow V^\pi(s) + \alpha(sample - V^\pi(s))$$

Example: TD Policy Evaluation

$$V^\pi(s) \leftarrow (1 - \alpha)V^\pi(s) + \alpha [R(s, \pi(s), s') + \gamma V^\pi(s')]$$

- | | |
|-----------------|-----------------|
| (1,1) up -1 | (1,1) up -1 |
| (1,2) up -1 | (1,2) up -1 |
| (1,2) up -1 | (1,3) right -1 |
| (1,3) right -1 | (2,3) right -1 |
| (2,3) right -1 | (3,3) right -1 |
| (3,3) right -1 | (3,2) up -1 |
| (3,2) up -1 | (4,2) exit -100 |
| (3,3) right -1 | (done) |
| (4,3) exit +100 | |
| (done) | |

Take $\gamma = 1, \alpha = 0.5$



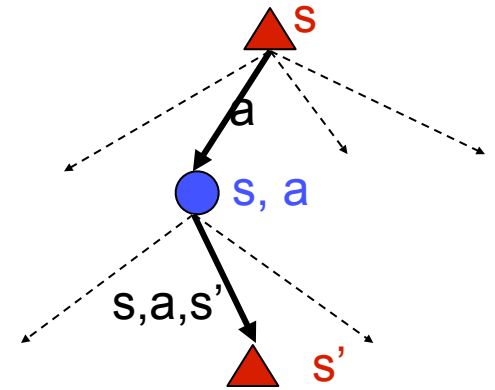
Problems with TD Value Learning

- TD value learning is model-free for policy evaluation (passive learning)
- However, if we want to turn our value estimates into a policy, we're sunk:

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

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

- Idea: learn Q-values directly
- Makes action selection model-free too!



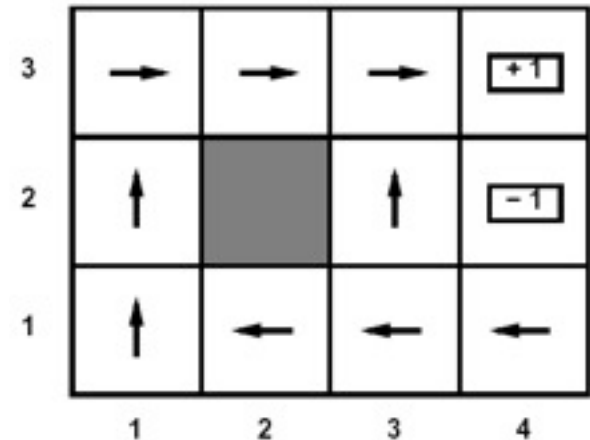
Active Learning

- Full reinforcement learning

- You don't know the transitions $T(s,a,s')$
- You don't know the rewards $R(s,a,s')$
- You can choose any actions you like
- **Goal: learn the optimal policy**
- ... what value iteration did!

- In this case:

- Learner makes choices!
- Fundamental tradeoff: exploration vs. exploitation
- This is NOT offline planning! You actually take actions in the world and find out what happens...



Detour: Q-Value Iteration

- Value iteration: find successive approx optimal values
 - Start with $V_0^*(s) = 0$
 - Given V_i^* , calculate the values for all states for depth $i+1$:

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_i(s')]$$

- But Q-values are more useful!
 - Start with $Q_0^*(s, a) = 0$
 - Given Q_i^* , calculate the q-values for all q-states for depth $i+1$:

$$Q_{i+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_i(s', a')]$$

Q-Learning Update

- Q-Learning: sample-based Q-value iteration

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

- Learn $Q^*(s, a)$ values

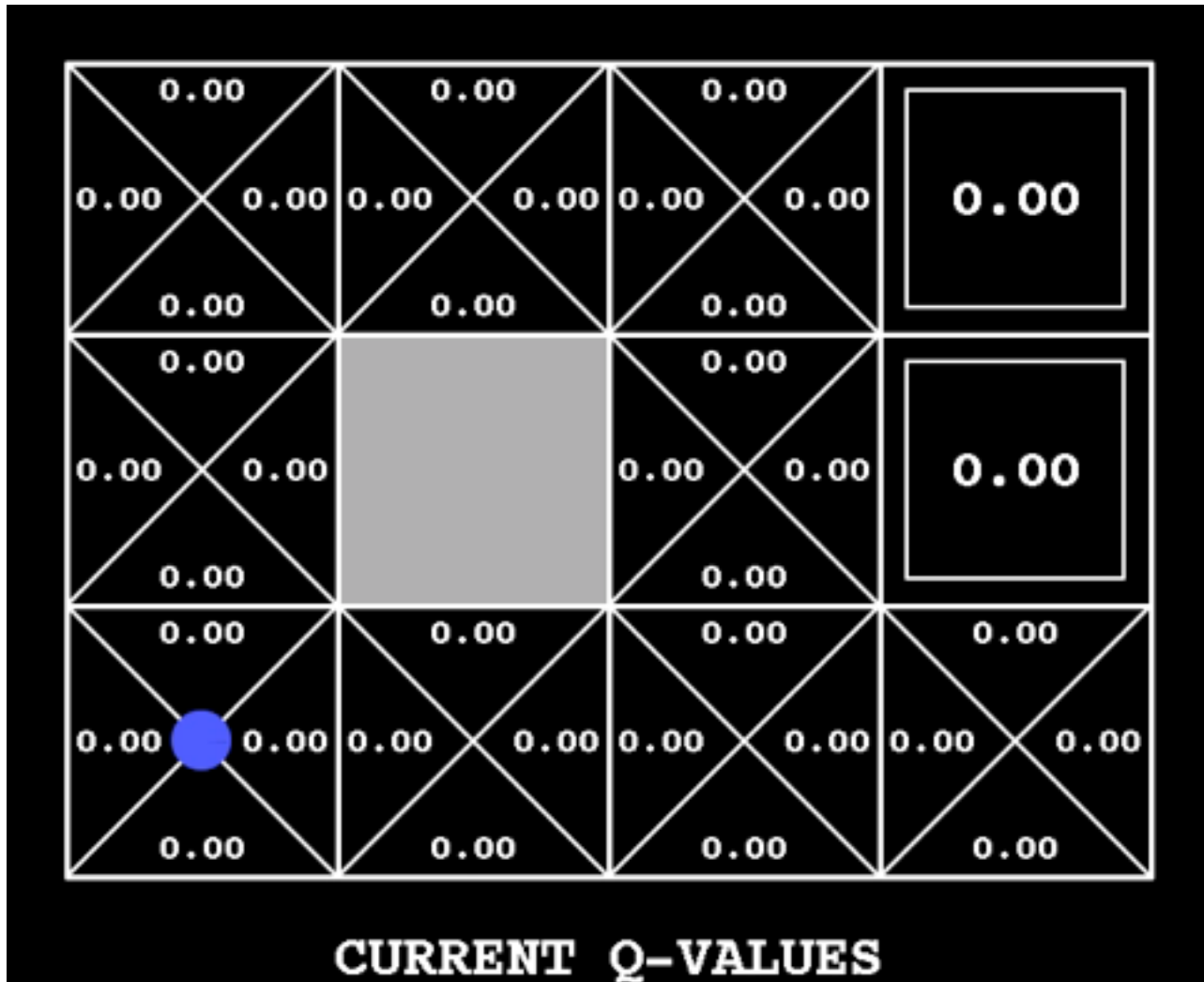
- Receive a sample (s, a, s', r)
- Consider your old estimate: $Q(s, a)$
- Consider your new sample estimate:

$$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

- Incorporate the new estimate into a running average:

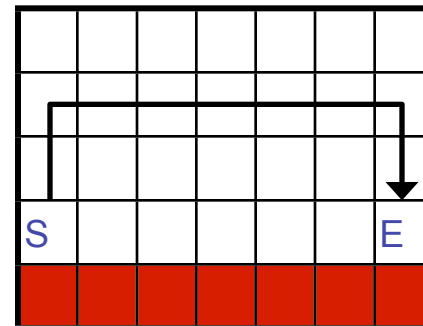
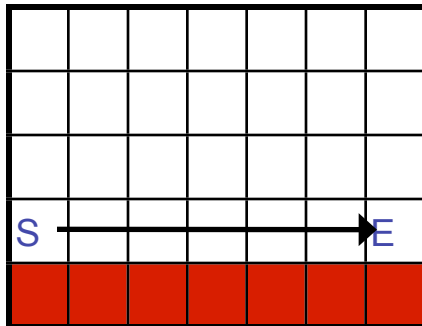
$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) [sample]$$

Q-Learning: Fixed Policy



Q-Learning Properties

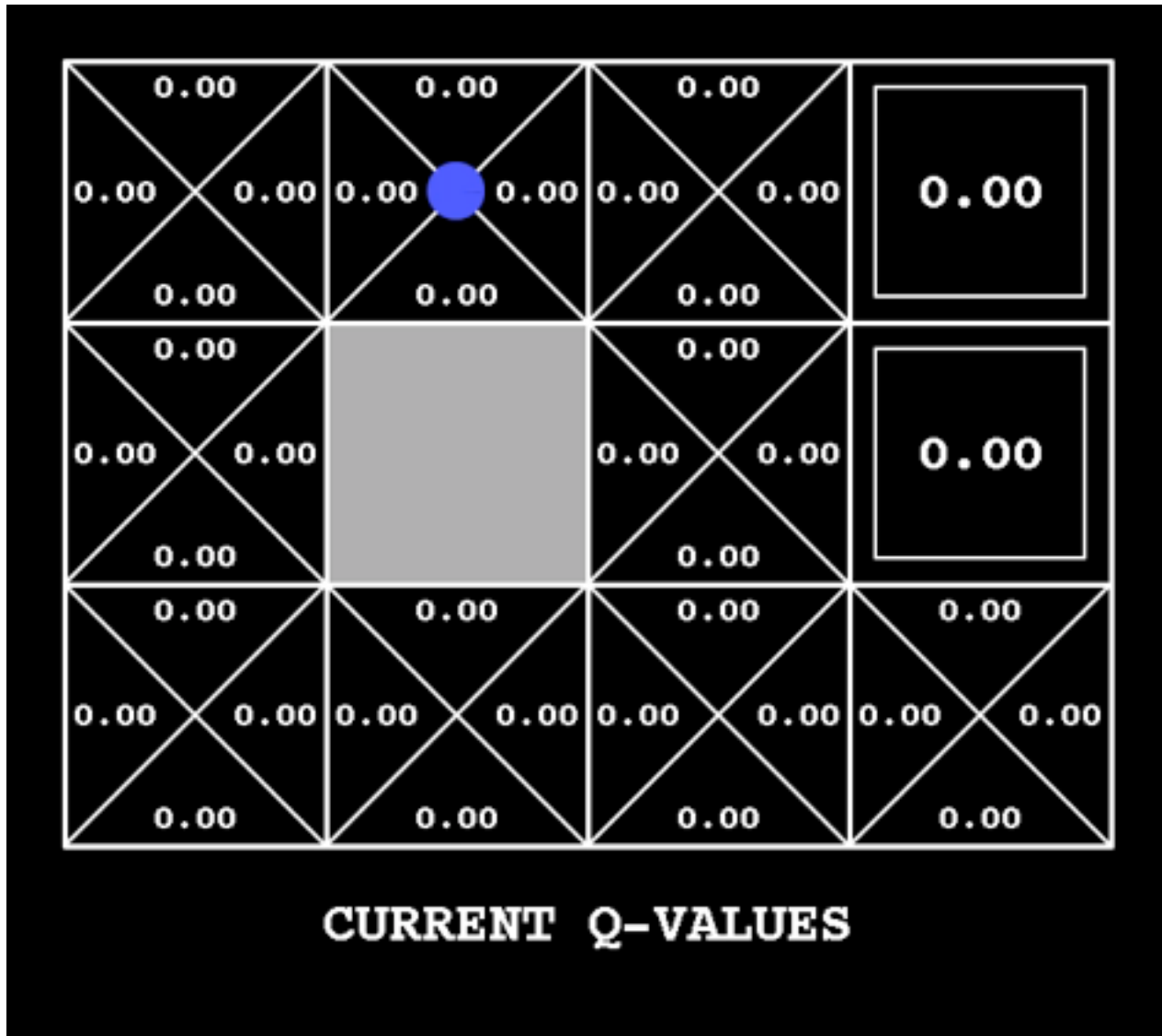
- Amazing result: Q-learning converges to optimal policy
 - If you explore enough
 - If you make the learning rate small enough
 - ... but not decrease it too quickly!
 - Not too sensitive to how you select actions (!)
- Neat property: off-policy learning
 - learn optimal policy without following it (some caveats)



Exploration / Exploitation

- Several schemes for action selection
 - Simplest: random actions (ϵ greedy)
 - Every time step, flip a coin
 - With probability ϵ , act randomly
 - With probability $1-\epsilon$, act according to current policy
 - Problems with random actions?
 - You do explore the space, but keep thrashing around once learning is done
 - One solution: lower ϵ over time
 - Another solution: exploration functions

Q-Learning: ϵ Greedy



Exploration Functions

- When to explore
 - Random actions: explore a fixed amount
 - Better idea: explore areas whose badness is not (yet) established
- Exploration function
 - Takes a value estimate and a count, and returns an optimistic utility, e.g. $f(u, n) = u + k/n$ (exact form not important)
 - Exploration policy $\pi(s)=$

$$\max_{a'} Q_i(s', a') \quad \text{vs.} \quad \max_{a'} f(Q_i(s', a'), N(s', a'))$$

Q-Learning Final Solution

- Q-learning produces tables of q-values:

