# CS 188: Artificial Intelligence
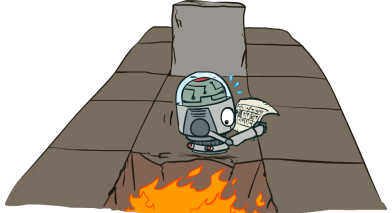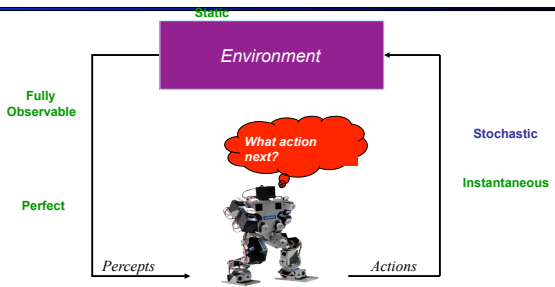## Markov Decision Processes II



Instructors: Dan Klein and Pieter Abbeel --- University of California, Berkeley

---

## Stochastic Planning: MDPs



Static

Environment

Fully Observable

Perfect

What action next?

Stochastic

Instantaneous

*Percepts*                    *Actions*
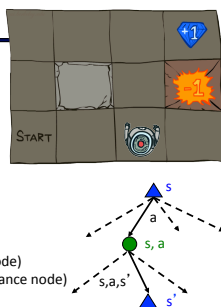
3

---

## 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 $\gamma$)
  - Start state $s_0$

- Quantities:
  - Policy = map of states to actions
  - Utility = sum of discounted rewards
  - Values = expected future utility from a state (max node)
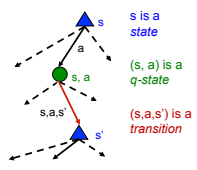  - Q-Values = expected future utility from a q-state (chance node)

s

a

s, a

s,a,s'

s'

---

## Solving MDPs



- Value Iteration
- ~~Real-Time Dynamic Programming~~
- Policy Iteration

- Reinforcement Learning

---

## 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

a

s, a

s,a,s'

s'

s is a *state*

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

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

[Demo: gridworld values (L9D1)]

---

## Gridworld Values V*



| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| 0.57 | | 0.57 | −1.00 |
| 0.49 | ◂ 0.43 | 0.48 | ◂ 0.28 |

VALUES AFTER 100 ITERATIONS

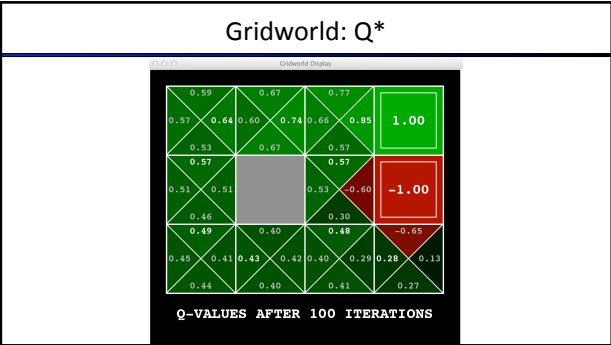## Gridworld: Q*



Q-VALUES AFTER 100 ITERATIONS

## 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

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

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

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

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



## 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 γ < 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



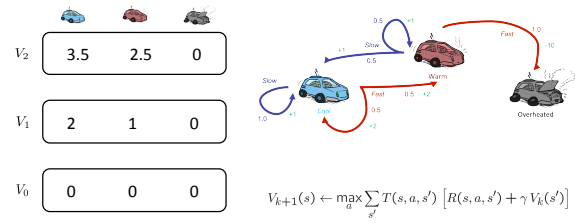$V_2(\ )$

[Demo – time-limited values (L8D6)]

## Time-Limited Values: Avoiding Redundant Computation

$V_4(\ )\ V_4(\ )\ V_4(\ )$

$V_3(\ )\ V_3(\ )\ V_3(\ )$

$V_2(\ )\ V_2(\ )\ V_2(\ )$

$V_1(\ )\ V_1(\ )\ V_1(\ )$

$V_0(\ )\ V_0(\ )\ V_0(\ )$

## Value Iteration



## Example: Value Iteration



| | | | |
|---|---|---|---|
| $V_2$ | 3.5 | 2.5 | 0 |
| $V_1$ | 2 | 1 | 0 |
| $V_0$ | 0 | 0 | 0 |

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

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

---

## Value Iteration

> Called a "Bellman Backup"
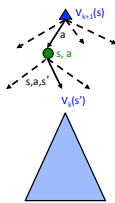
- **Start with $V_0(s) = 0$:** *no time steps left means an expected reward sum of zero*
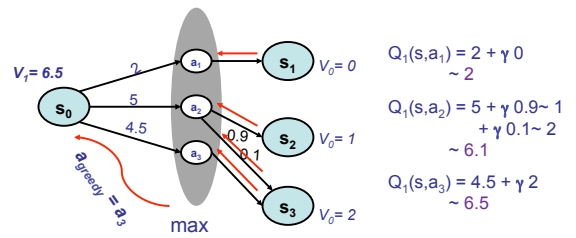- **Given vector of $V_k(s)$ values, do one ply of expectimax from each state:**

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

- **Repeat until convergence**
  *(trust me, it does)*



---

## Example: Bellman Backup



$V_1 = 6.5$

$V_0 = 0$
$V_0 = 1$
$V_0 = 2$

$a_{greedy} = a_3$

max

$Q_1(s,a_1) = 2 + \gamma \, 0$
$\sim 2$

$Q_1(s,a_2) = 5 + \gamma \, 0.9 \sim 1$
$\qquad + \gamma \, 0.1 \sim 2$
$\qquad \sim 6.1$

$Q_1(s,a_3) = 4.5 + \gamma \, 2$
$\qquad \sim 6.5$

---

## k=0



VALUES AFTER 0 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



| | | | |
|---|---|---|---|
| 0.00 ▲ | 0.00 ▸ | 0.72 ▸ | 1.00 |
| 0.00 ▲ | | 0.00 ▲ | −1.00 |
| 0.00 ▲ | 0.00 ▲ | 0.00 ▲ | 0.00 ▾ |

VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=3



| | | | |
|---|---|---|---|
| 0.00 ▸ | 0.52 ▸ | 0.78 ▸ | 1.00 |
| 0.00 ▲ | | 0.43 | −1.00 |
| 0.00 ▲ | 0.00 ▲ | 0.00 ▲ | 0.00 ▾ |

VALUES AFTER 3 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=4



| | | | |
|---|---|---|---|
| 0.37 ▸ | 0.66 ▸ | 0.83 ▸ | 1.00 |
| 0.00 ▲ | | 0.51 ▲ | −1.00 |
| 0.00 ▲ | 0.00 ▸ | 0.31 ▲ | ◂ 0.00 |

VALUES AFTER 4 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0
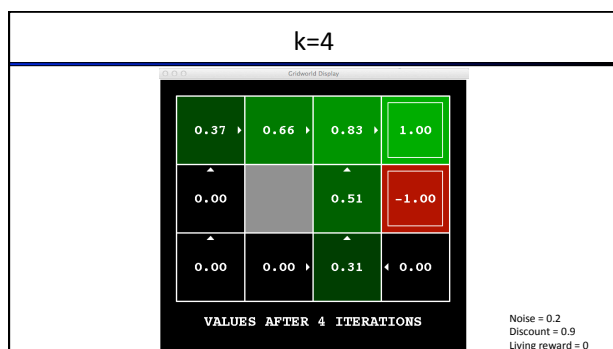
## k=5



| | | | |
|---|---|---|---|
| 0.51 ▸ | 0.72 ▸ | 0.84 ▸ | 1.00 |
| 0.27 ▲ | | 0.55 ▲ | −1.00 |
| 0.00 ▲ | 0.22 ▸ | 0.37 ▲ | ◂ 0.13 |

VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=6



| | | | |
|---|---|---|---|
| 0.59 ▸ | 0.73 ▸ | 0.85 ▸ | 1.00 |
| 0.41 ▲ | | 0.57 ▲ | −1.00 |
| 0.21 ▲ | 0.31 ▸ | 0.43 ▲ | ◂ 0.19 |

VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=7



| | | | |
|---|---|---|---|
| 0.62 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| 0.50 ▲ | | 0.57 ▲ | −1.00 |
| 0.34 ▲ | 0.36 ▸ | 0.45 ▲ | ◂ 0.24 |

VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=8



| | | | |
|---|---|---|---|
| 0.63 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.53 | | 0.57 | −1.00 |
| ▲ 0.42 | 0.39 ▸ | ▲ 0.46 | ◂ 0.26 |

VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=9



| | | | |
|---|---|---|---|
| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.55 | | 0.57 | −1.00 |
| ▲ 0.46 | 0.40 ▸ | ▲ 0.47 | ◂ 0.27 |

VALUES AFTER 9 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=10



| | | | |
|---|---|---|---|
| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.56 | | 0.57 | −1.00 |
| ▲ 0.48 | ◂ 0.41 | ▲ 0.47 | ◂ 0.27 |

VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=11



| | | | |
|---|---|---|---|
| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.56 | | 0.57 | −1.00 |
| ▲ 0.48 | ◂ 0.42 | ▲ 0.47 | ◂ 0.27 |

VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=12



| | | | |
|---|---|---|---|
| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.57 | | 0.57 | −1.00 |
| ▲ 0.49 | ◂ 0.42 | ▲ 0.47 | ◂ 0.28 |

VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

## k=100



| | | | |
|---|---|---|---|
| 0.64 ▸ | 0.74 ▸ | 0.85 ▸ | 1.00 |
| ▲ 0.57 | | 0.57 | −1.00 |
| ▲ 0.49 | ◂ 0.43 | ▲ 0.48 | ◂ 0.28 |

VALUES AFTER 100 ITERATIONS

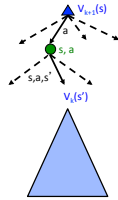Noise = 0.2
Discount = 0.9
Living reward = 0

## Value Iteration

- **Start with $V_0(s) = 0$:**

- **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**

- **Complexity of each iteration: O(S²A)**
- **Number of iterations: poly(|S|, |A|, 1/(1-g))**
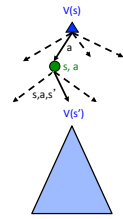
- **Theorem: will converge to unique optimal values**

---

## Value Iteration
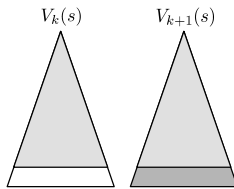
- Bellman equations characterize the optimal values:

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

- Value iteration computes them:

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

- Value iteration is just a fixed point solution method
  - ... though the $V_k$ vectors are also interpretable as time-limited values
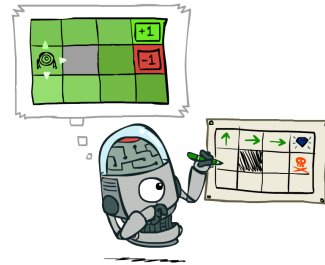
---

## Convergence*

- How do we know the $V_k$ vectors will converge?

- Case 1: If the tree has maximum depth M, then $V_M$ holds the actual untruncated values

- Case 2: If the discount is less than 1
  - 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 max difference happens if big reward at k+1 level
  - That last layer is at best all $R_{MAX}$
  - 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

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

---

## Policy Extraction



---

## Computing Actions from Values

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

- How should we act?
  - 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

| 0.95 ▸ | 0.96 ▸ | 0.98 ▸ | 1.00 |
| 0.94 | | ◂ 0.89 | −1.00 |
| 0.92 | ◂ 0.91 | ◂ 0.90 | 0.80 |

---

## Computing Actions from Q-Values

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

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

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

- Important lesson: actions are easier to select from q-values than values!
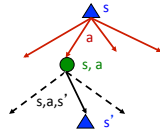
## Problems with Value Iteration

- Value iteration repeats the Bellman updates:

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

- Problem 1: It's slow – $O(S^2A)$ per iteration

- Problem 2: The "max" at each state rarely changes

- Problem 3: The policy often converges long before the values

[Demo: value iteration (L9D2)]

## VI → Asynchronous VI

- Is it essential to back up *all* states in each iteration?
  - No!

- States may be backed up
  - many times or not at all
  - in any order

- As long as no state gets starved...
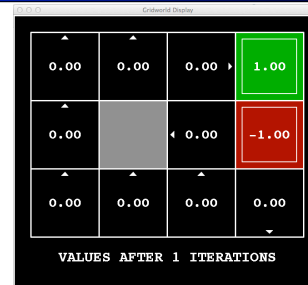  - convergence properties still hold!!

44

## Prioritization of Bellman Backups

- Are all backups equally important?

- Can we avoid some backups?

- Can we schedule the backups more appropriately?

45

## k=1



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
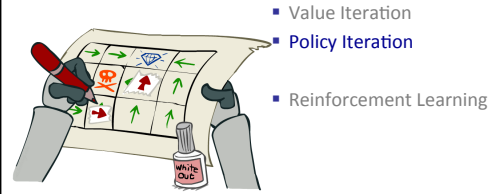
## k=3



VALUES AFTER 3 ITERATIONS

Noise = 0.2
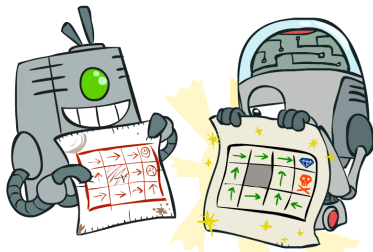Discount = 0.9
Living reward = 0

## Asynch VI: Prioritized Sweeping

- Why backup a state if values of successors same?
- Prefer backing a state
  - whose successors had most change

- Priority Queue of (state, expected change in value)
- Backup in the order of priority
- After backing a state update priority queue
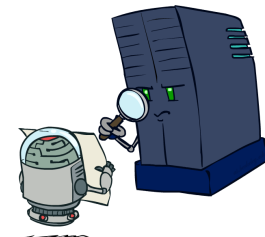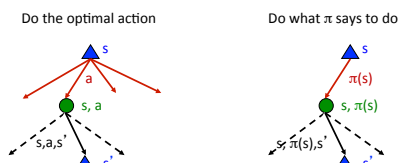  - for all predecessors

## Solving MDPs



- Value Iteration
- Policy Iteration
- Reinforcement Learning

## Policy Methods



## Policy Evaluation



## Fixed Policies

Do the optimal action          Do what π says to do



s                              s
a                              π(s)
s, a                           s, π(s)
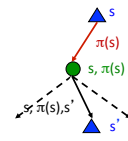s,a,s'                         s, π(s),s'
s'                             s'

- Expectimax trees max over all actions to compute the optimal values

- If we fixed some policy π(s), then the tree would be simpler – only one action per state
  - … though the tree's value would depend on which policy we fixed

## Utilities for a Fixed Policy

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

- Define the utility of a state s, under a fixed policy π:
  Vπ(s) = expected total discounted rewards starting in s and following π

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

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

s
π(s)
s, π(s)
s, π(s),s'
s'

8

## Example: Policy Evaluation

Always Go Right            Always Go Forward



---

## Example: Policy Evaluation

Always Go Right            Always Go Forward



---

## Policy Evaluation

- How do we calculate the V's for a fixed policy $\pi$?

- Idea 1: Turn recursive Bellman equations into updates (like value iteration)

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

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

- Efficiency: O(S$^2$) per iteration

- Idea 2: Without the maxes, the Bellman equations are just a linear system
  - Solve with Matlab (or your favorite linear system solver)



---

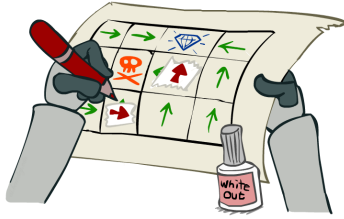## Policy Extraction



---

## Computing Actions from Values

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

- How should we act?
  - 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!

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

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

## Policy Iteration



## Policy Iteration

- Alternative approach for optimal values:
  - Step 1: Policy evaluation: calculate utilities for some fixed policy (not optimal utilities!) until convergence
  - Step 2: Policy improvement: update policy using one-step look-ahead with resulting converged (but not optimal!) utilities as future values
  - Repeat steps until policy converges

- This is policy iteration
  - It's still optimal!
  - Can converge (much) faster under some conditions

## Policy Iteration

- Evaluation: For fixed current policy π, find values with policy evaluation:
  - Iterate until values converge:

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

- Improvement: For fixed values, get a better policy using policy extraction
  - One-step look-ahead:

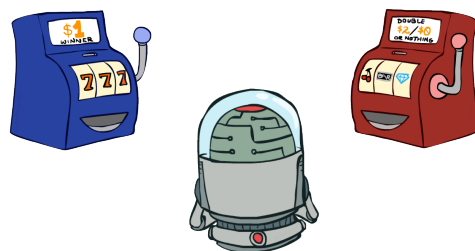$$\pi_{i+1}(s) = \arg\max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^{\pi_i}(s') \right]$$
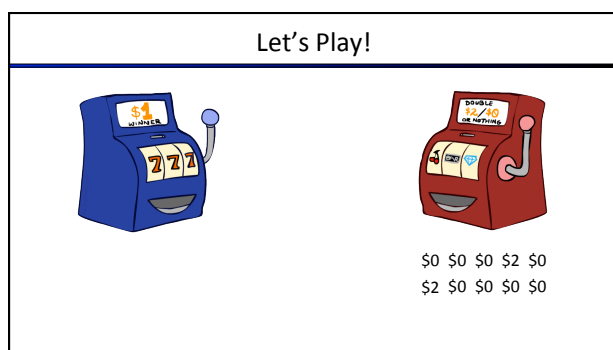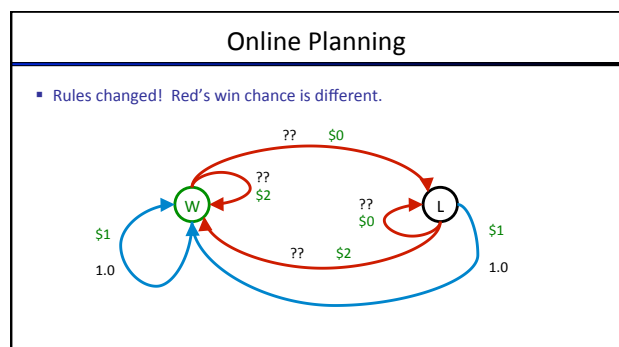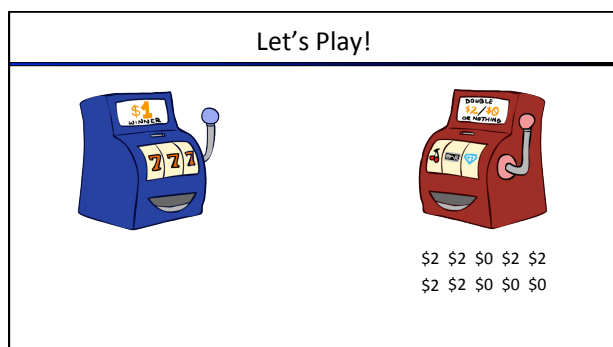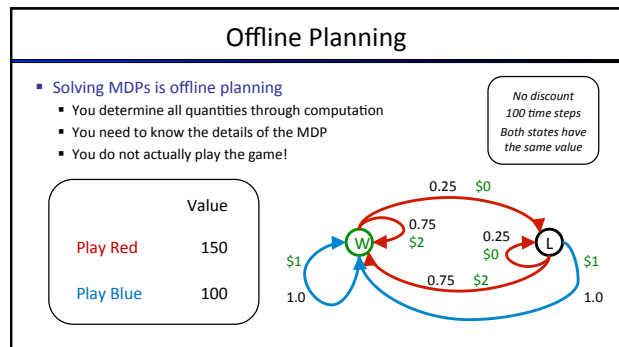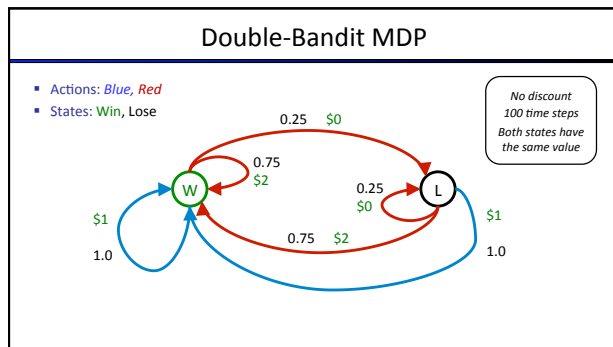
## Comparison

- Both value iteration and policy iteration compute the same thing (all optimal values)

- In value iteration:
  - Every iteration updates both the values and (implicitly) the policy
  - We don't track the policy, but taking the max over actions implicitly recomputes it

- In policy iteration:
  - We do several passes that update utilities with fixed policy (each pass is fast because we consider only one action, not all of them)
  - After the policy is evaluated, a new policy is chosen (slow like a value iteration pass)
  - The new policy will be better (or we're done)

- Both are dynamic programs for solving MDPs

## Summary: MDP Algorithms

- So you want to....
  - Compute optimal values: use value iteration or policy iteration
  - Compute values for a particular policy: use policy evaluation
  - Turn your values into a policy: use policy extraction (one-step lookahead)

- These all look the same!
  - They basically are – they are all variations of Bellman updates
  - They all use one-step lookahead expectimax fragments
  - They differ only in whether we plug in a fixed policy or max over actions

## Double Bandits

## Double-Bandit MDP

- Actions: *Blue*, *Red*
- States: Win, Lose

*No discount*
*100 time steps*
*Both states have*
*the same value*



0.25  $0
0.75  $2
0.25  $0
0.75  $2
$1
1.0
$1
1.0

---

## Offline Planning

- Solving MDPs is offline planning
  - You determine all quantities through computation
  - You need to know the details of the MDP
  - You do not actually play the game!

*No discount*
*100 time steps*
*Both states have*
*the same value*

|  | Value |
|---|---|
| Play Red | 150 |
| Play Blue | 100 |



0.25  $0
0.75  $2
0.25  $0
0.75  $2
$1
1.0
$1
1.0

---

## Let's Play!



$2  $2  $0  $2  $2
$2  $2  $0  $0  $0

---

## Online Planning

- Rules changed!  Red's win chance is different.



??  $0
??  $2
??  $0
??  $2
$1
1.0
$1
1.0

---

## Let's Play!



$0  $0  $0  $2  $0
$2  $0  $0  $0  $0

---

## What Just Happened?

- That wasn't planning, it was learning!
  - Specifically, reinforcement learning
  - There was an MDP, but you couldn't solve it with just computation
  - You needed to actually act to figure it out

- Important ideas in reinforcement learning that came up
  - Exploration: you have to try unknown actions to get information
  - Exploitation: eventually, you have to use what you know
  - Regret: even if you learn intelligently, you make mistakes
  - Sampling: because of chance, you have to try things repeatedly
  - Difficulty: learning can be much harder than solving a known MDP

Next Time: Reinforcement Learning!