# Reinforcement Learning:
# A brief introduction

**Mihaela van der Schaar**

# Outline

- **Optimal Decisions & Optimal Forecasts –**

- **Markov Decision Processes (MDPs)**

  - States, actions, rewards and value functions

  - Dynamic Programming

  - Bellman equations

    - Value Iteration

    - Policy Iteration

  - Illustrative examples

  - Reinforcement learning

# Markov decision process (MDP)

- Discrete-time stochastic *control* process

- Extension of Markov chains

- Differences:
  - Addition of actions (choice)
  - Addition of rewards (goal)

- If the actions are fixed and there are no rewards, an MDP reduces to a Markov chain

# Discrete MDP model

- Time $t$ is discrete
- State space $S$
- Set of actions $A$
- Reward function $R(s, a)$
- Transition model $p(s'|s, a)$

Cost function can be used instead!

$$C(s, a)$$

(Cost = -Reward)

The Markov property entails that the next state $s_{t+1}$ only depends on the previous state $s_t$ and action $a_t$:

$$p(s_{t+1}|s_t, s_{t-1}, \ldots, s_0, a_t, a_{t-1}, \ldots, a_0) = p(s_{t+1}|s_t, a_t)$$

# Rewards and optimality criterion

Agent should maximize

(minimize)

$$E\left[\sum_{t=0}^{h} \gamma^t R_t\right], \quad \left(E\left[\sum_{t=0}^{h} \gamma^t C_t\right]\right) \qquad (2)$$

where

- $h$ is the planning horizon, can be finite or $\infty$
- $\gamma$ is a discount rate, $0 \leq \gamma < 1$    Myopic vs. foresighted decisions

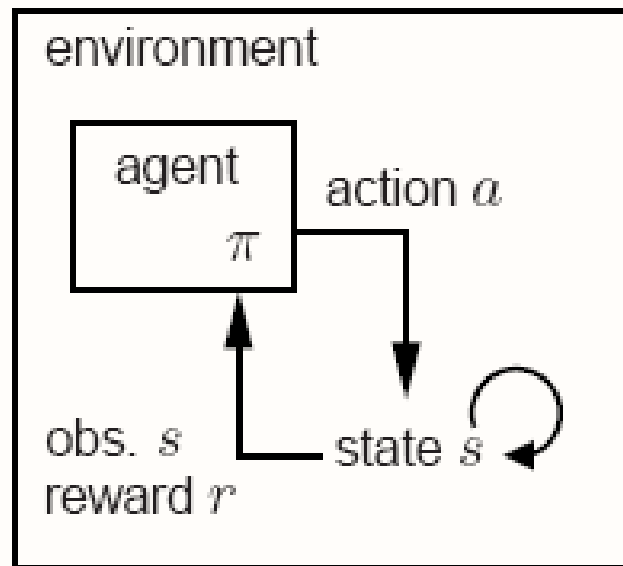Reward hypothesis (Sutton and Barto, 1998):
All goals and purposes can be formulated as the maximization of the cumulative sum of a received scalar signal (reward).

Examples?

# Solution to an MDP = Policy $\pi$

- Gives the action to take from a given state.
- Does not depend on history (Why?)
- Goal: Find a policy that maximizes the cumulative discounted sum of rewards $\left( E\left[ \sum_{t=0}^{h} \gamma^t R_t \right] \right)$

# Brief summary of concepts

- The *agent* and its *environment* interact over a sequence of discrete time steps.

- The specification of their interface defines a particular task:

  - the *actions* are the choices made by the agent;

  - the *states* are the basis for making the choices;

  - the *rewards* are the basis for evaluating the choices.

- A *policy* is a stochastic rule by which the agent selects actions as a function of states.

- The agent's objective is to maximize the amount of reward it receives over time.

# Policies and value

An agent acts according to its policy

$$\pi : S \to A. \tag{3}$$

A common way to characterize a policy is by its $\boxed{\text{value function:}}$

**Quantifies how good it is to be in a state**

$$V^{\pi}(s) = R(s, \pi(s)) + E\left[ \sum_{t=1}^{\infty} \gamma^t R(s_t, \pi(s_t)) \right]. \tag{4}$$

**State-value function**

The expectation operator averages over the stochastic transition model, which leads to the following recursion:

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s'|s, \pi(s)) V^{\pi}(s'). \tag{5}$$

# Policies and value

Extracting a policy $\pi$ from a value function $V$ is easy:

$$\pi(s) = \arg\max_{a \in A} \left[ R(s,a) + \gamma \sum_{s' \in S} p(s'|s,a) V(s') \right]. \quad (6)$$

Bellman (1957) equation:

$$V^*(s) = \max_\pi V^\pi(s)$$

$$V^*(s) = \max_{a \in A} \left[ R(s,a) + \gamma \sum_{s' \in S} p(s'|s,a) V^*(s') \right], \quad (7)$$

**Optimal state-value function**

**Known**

Solving this (nonlinear) system of equations for each state $s$ yields the optimal value function, and an optimal policy $\pi^*$. However, due to the nonlinear $\max$ operator solving the system for each state simultaneously is not efficient for large MDPs

# Brief summary of concepts

- A policy's *value function* assigns to each state the expected return from that state given that the agent uses the policy.

- The *optimal value function* assigns to each state the largest expected return achievable by any policy.

- A policy whose value functions are optimal is an *optimal policy*.

- The optimal value functions for states are unique for a given MDP, but there can be many optimal policies.

- The *Bellman optimality equations* are special consistency conditions that the optimal value functions must satisfy and that can, in principle, be solved to obtain the optimal value functions, from which an optimal policy can be determined with relative ease.

# How do we compute the optimal state-value function and the optimal policy?

- Problem:
  - Given:
    - Transition probability function: $p\left( s' \mid s, a \right)$
    - Reward function: $R(s, a)$
  - Determine:
    - Optimal state-value function: $V^*$
    - Optimal policy: $\pi^*$

Solution:

Dynamic Programming

# Dynamic Programming

- The term dynamic programming (DP) refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment
  - Can be used to solve Markov decision processes.
  - Use value functions to organize and structure the search for good policies.
  - Turn Bellman equations into update policies.

# DP methods

- *Policy evaluation* refers to the (typically) iterative computation of the value functions for a given policy.

- *Policy improvement* refers to the computation of an improved policy given the value function for that policy.

- Putting these two computations together, we obtain *policy iteration* and *value iteration*, the two most popular DP methods.

    - Either of these can be used to reliably compute optimal policies and value functions for finite MDPs given complete knowledge of the MDP.

# Policy evaluation

Input $\pi$, the policy to be evaluated

Initialize $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat

    $\Delta \leftarrow 0$

    For each $s \in \mathcal{S}$:

        $v \leftarrow V(s)$

        $V(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' \mid s, \pi(s)) V(s')$

        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number)

Output $V \approx V^\pi$

# Policy improvement

- When should we change the policy?
  - Let $\pi$ be the current policy
  - Let $\pi' = \pi$ except in state $s$ where we let $\pi'(s) = \alpha$
  - If we pick a new action $\alpha$ to take in state $s$, and $V(\pi') \geq V(\pi)$, then picking $\alpha$ from state $s$ is a better policy overall.

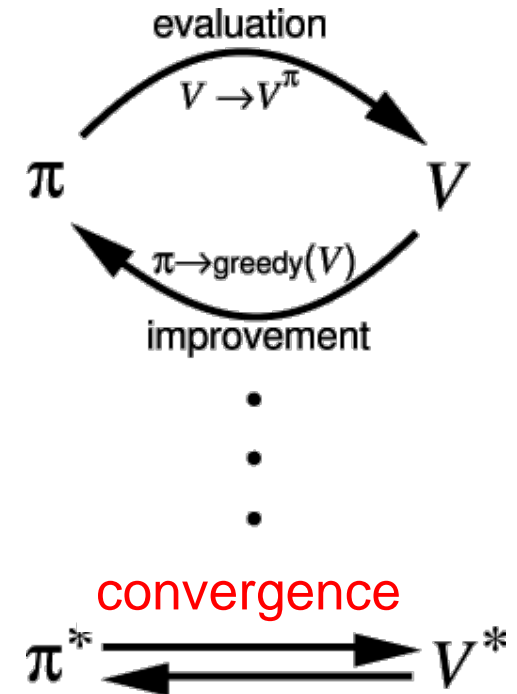- The above follows from the *policy improvement theorem*

# Iterative DP algorithms

- All iterative algorithms for computing the optimal state-value function and optimal policy have two basic steps:
  - 1. Policy Improvement
  - 2. Policy Evaluation

1. $$\pi(s) := \arg\max_{a \in A} \left\{ R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right\}$$

2. $$V(s) := R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' \mid s, \pi(s)) V(s')$$

Value Function

evaluation

$V \to V^\pi$

$\pi$      $V$

$\pi \to \text{greedy}(V)$

improvement

convergence

$\pi^* \rightleftarrows V^*$

Next, several algorithms that apply these two basic steps
(in different orders) for computing optimal value function and optimal policy

# Value iteration

**can be used to compute optimal policies and value functions**

Value iteration: successive approximation technique.
The optimal value function $V_0^*$

$$V_0^*(s) = \max_{a \in A} R(s, a). \tag{8}$$

arbitrary initial value, e.g. optimal myopic value

In order to consider one step deeper into the future, i.e., to compute $V_{n+1}^*$ from $V_n^*$

$$V_{n+1}^*(s) = \max_{a \in A} \left[ R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V_n^*(s') \right], \tag{9}$$

known

which is known as a Bellman backup $H$, allowing us to write (9) as

$$V_{n+1}^* = HV_n^*. \tag{10}$$

Note that policy improvement and evaluation in (9) are combined (simultaneous)

# Value iteration

Initialize $V$ arbitrarily, e.g., $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat
$\quad \Delta \leftarrow 0$
$\quad$ For each $s \in \mathcal{S}$:
$\qquad v \leftarrow V(s)$
$\qquad V(s) \leftarrow \max_{a \in A} \left\{ R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s,a) V(s') \right\}$
$\qquad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, $\pi$, such that
$\quad \pi(s) \leftarrow \arg\max_{a \in A} \left\{ R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s,a) V(s') \right\}$

# Value iteration discussion

- As $n \to \infty$, value iteration converges.

- Value iteration has converged when the largest update $\Delta$ in an iteration is below a certain threshold $\theta$

- Exhaustive sweeps are not required for convergence: arbitrary states can be backed up in arbitrary order, provided that in the limit all states are visited infinitely often (Bertsekas and Tsitsiklis, 1989).

- This can be exploited by backing up the most promising states first, known as prioritized sweeping (Moore and Atkeson, 1993; Peng and Williams, 1993).

# Value iteration example: Grid World

- States: (x,y)-coordinate in a 10 x 10 grid
- Actions: Up, down, left, and right
- Dynamics:
  - Move one-step in direction specified by action with prob 0.7
  - Move one-step in any one of the other directions with prob 0.1
  - Cannot move outside of the grid (i.e. end up in the same state)
  - Agent is flung randomly to corner of grid after entering a Goal or Penalty state

- Rewards:
  - Attempted move outside of grid leads to reward of -1
  - Goal/Penalty states in grid: +10@(9,8), +3@(8,3), -5@(4,5), -10@(4,8)
  - No reward in any other states
- Discount factor: 0.9

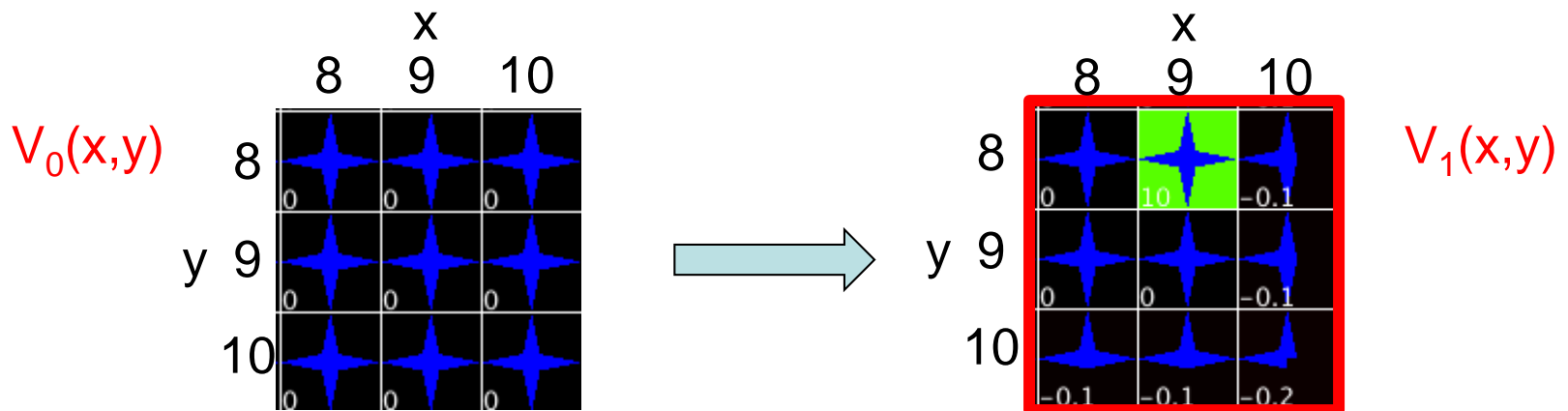# Value iteration example: Grid World



+3

-5

+10

-10

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_1(x,y)$ (Iteration 1)

How are these calculated?

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World

$$V_1^*\left(9,8\right) = \max_a \left[ R\left(9,8,a\right) + \gamma \sum_{(x',y')} p\left(x',y' \mid 9,8,a\right) V_0^*\left(x',y'\right)\right]$$

$$= \max_a R\left(9,8,a\right)$$

$$= \max_a \left\{ \underbrace{10}_{\text{up}}, \underbrace{10}_{\text{down}}, \underbrace{10}_{\text{left}}, \underbrace{10}_{\text{right}} \right\} = 10$$



$V_0(x,y)$
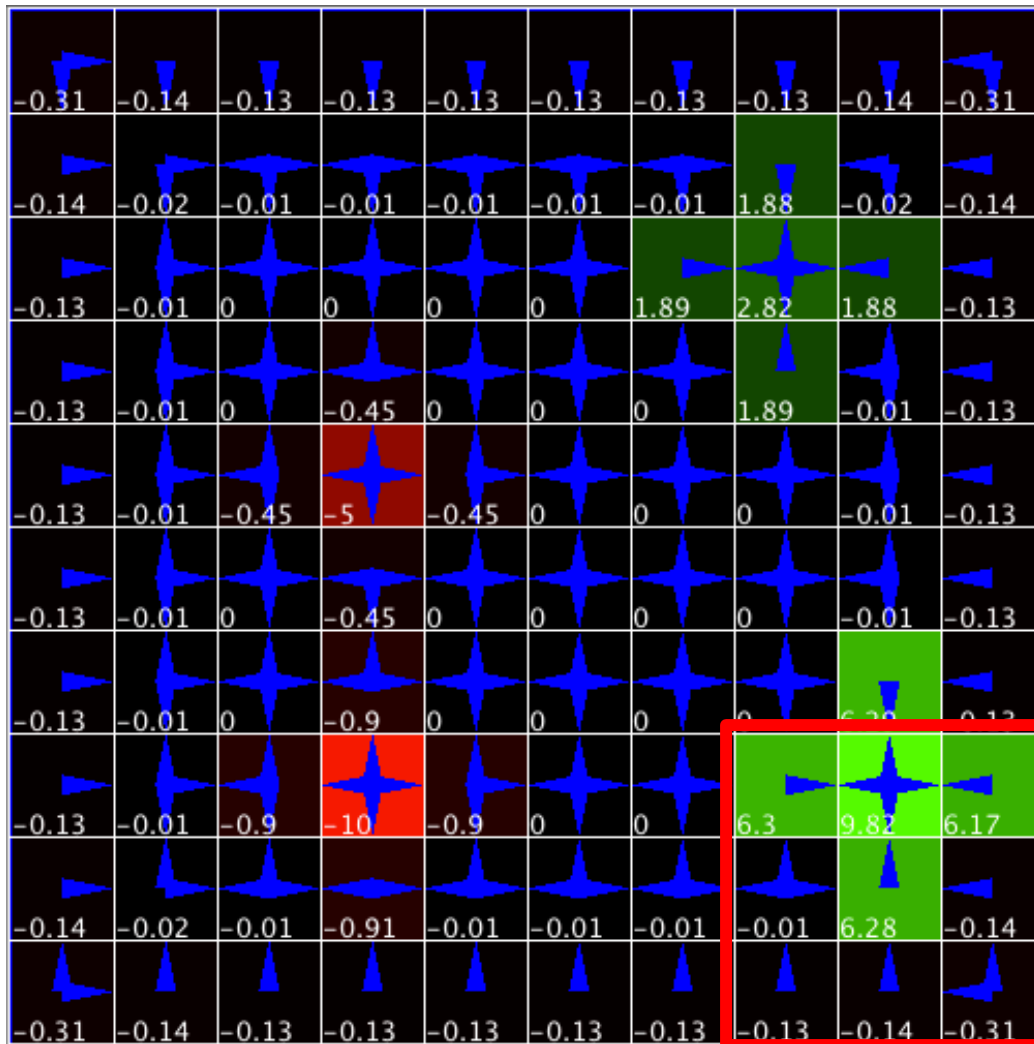
$V_1(x,y)$

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

23

# Value iteration example: Grid World

$$V_1^*\left(10,10\right) = \max_a \left[ R\left(10,10,a\right) + \gamma \sum_{\left(x',y'\right)} p\left(x',y' \mid 10,10,a\right)V_0^*\left(x',y'\right) \right]$$

$$= \max_a R\left(10,10,a\right)$$

$$= \max_a \left\{ \underbrace{2 \cdot 0.1\left(-1\right)}_{\text{up}}, \; \underbrace{\left(0.7 + 0.1\right)\left(-1\right)}_{\text{down}}, \; \underbrace{2 \cdot 0.1\left(-1\right)}_{\text{left}}, \; \underbrace{\left(0.7 + 0.1\right)\left(-1\right)}_{\text{right}} \right\} = -0.2$$



$V_0(x,y)$

$V_1(x,y)$

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_2(x,y)$ (Iteration 2)

How are these calculated?

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World

$$V_2^*\left(9,8\right) = \max_a \left[ R\left(9,8,a\right) + \gamma \sum_{(x',y')} p\left(x',y' \mid 9,8,a\right) V_1^*\left(x',y'\right)\right]$$

$$= 10 + 0.9 \times 0.25\left[V\left(1,1\right) + V(1,10) + V(10,1) + V(10,10)\right]$$

$$= 10 + 0.9 \times -0.2 = 9.82$$



10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World

$$V_2^*\left(10,10\right) = \max_a\left[R\left(10,10,a\right) + \gamma\sum_{(x',y')} p\left(x',y' \mid 10,10,a\right)V_1^*\left(x',y'\right)\right]$$

$a = \text{up}$      $-0.2 + 0.9\left[0.7\left(-0.1\right) + 0.1\left(-0.1\right) + 2{\cdot}0.1\left(-0.2\right)\right] = -0.31$

$a = \text{down}$      $-0.8 + 0.9\left[\left(0.7 + 0.1\right)\left(-0.2\right) + 2{\cdot}0.1\left(-0.1\right)\right] = -0.96$

$V_1(x,y)$

$V_2(10,10)$



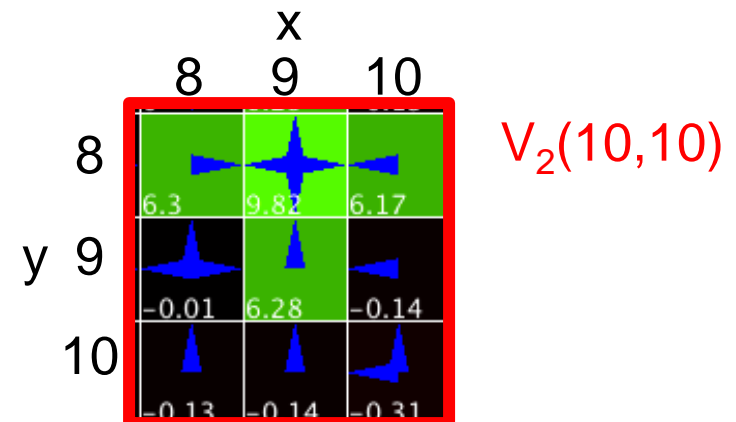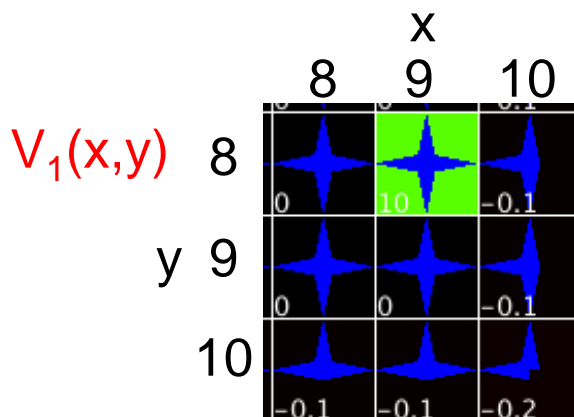10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World

$$V_2^*(10,10) = \max_a \left[ R(10,10,a) + \gamma \sum_{(x',y')} p(x',y' \mid 10,10,a) V_1^*(x',y') \right]$$

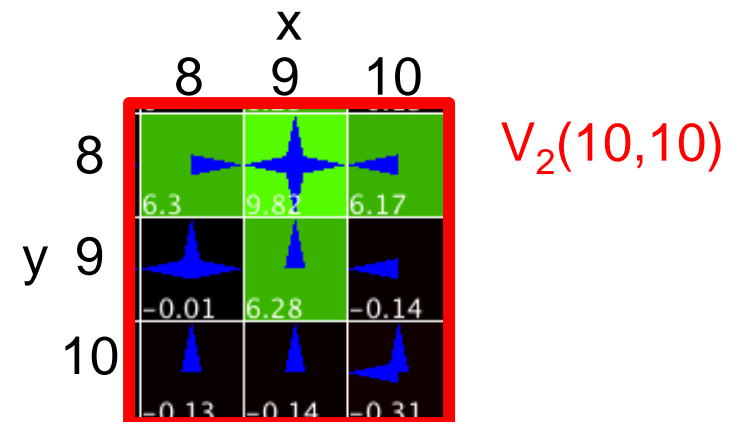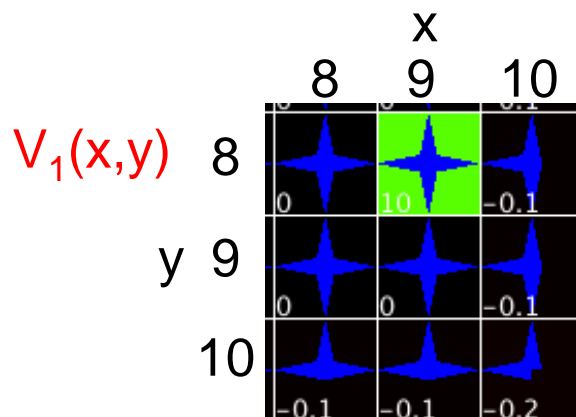$a = \text{left}$    $-0.2 + 0.9\left[0.7(-0.1) + 0.1(-0.1) + 2{\cdot}0.1(-0.2)\right] = -0.31$

$a = \text{right}$    $-0.8 + 0.9\left[(0.7 + 0.1)(-0.2) + 2{\cdot}0.1(-0.1)\right] = -0.96$



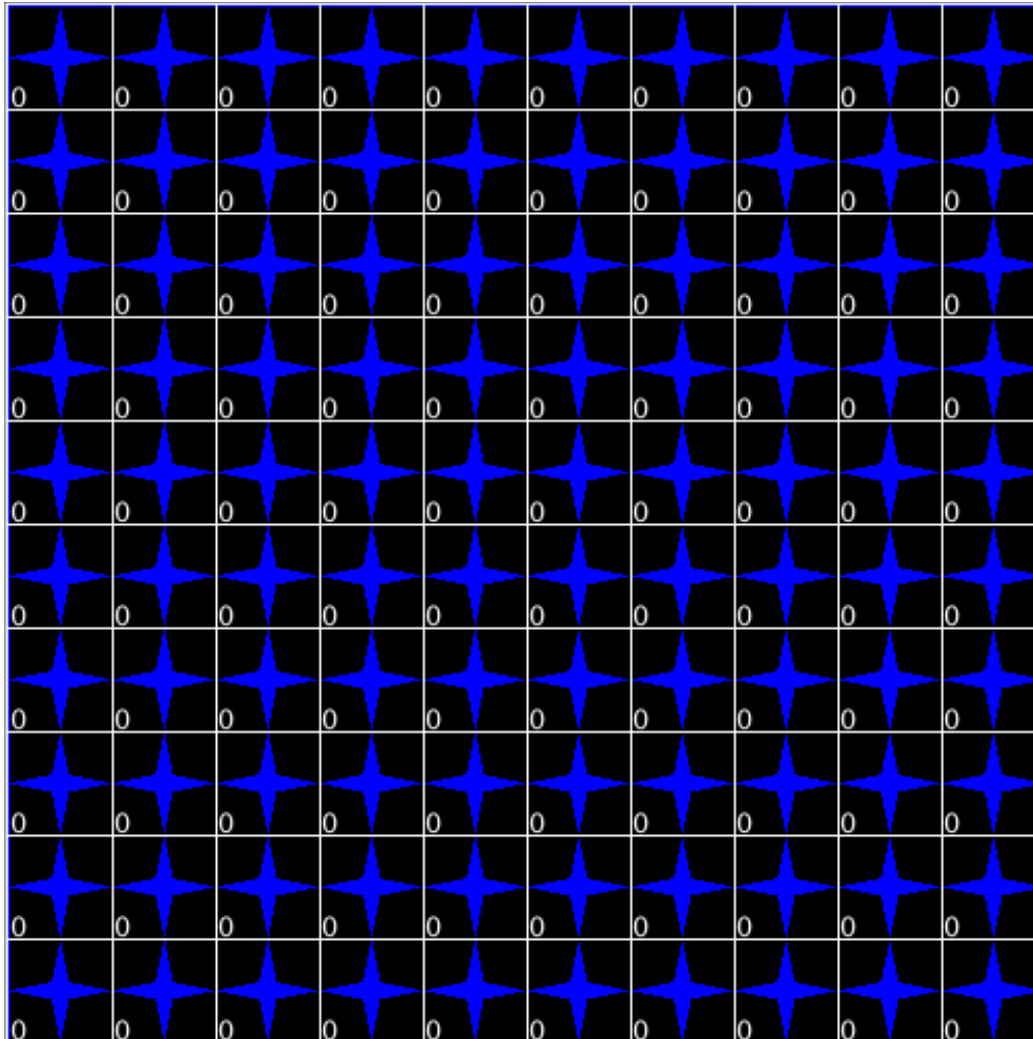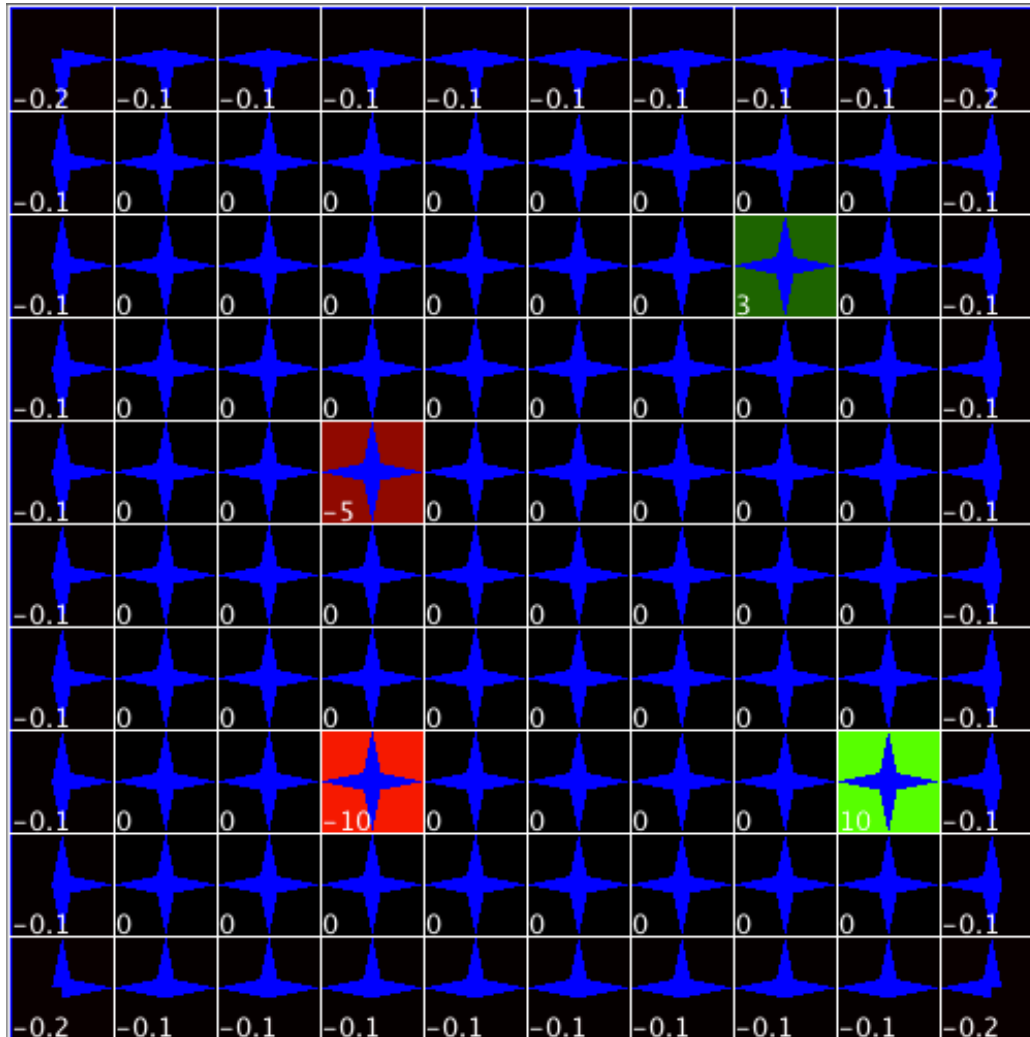10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World

$$V_2^*\left(10,10\right) = \max_a\left[R\left(10,10,a\right) + \gamma\sum_{(x',y')} p\left(x',y' \mid 10,10,a\right)V_1^*\left(x',y'\right)\right]$$

$$= -0.31$$



$V_1(x,y)$

$V_2(10,10)$

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_0(x,y)$ (Iteration 0)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_1(x,y)$ (Iteration 1)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_2(x,y)$ (Iteration 2)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_2(x,y)$ (Iteration 3)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_{10}(x,y)$ (Iteration 10)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_{20}(x,y)$ (Iteration 20)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Value iteration example: Grid World



$V_{50}(x,y)$ (Iteration 50)

10 x 10 Grid World Java Applet: http://www.cs.ubc.ca/~poole/demos/mdp/vi.html

# Policy iteration

- Iteratively evaluate $V^\pi$ and improve the policy $\pi$

- A finite MDP has a finite number of policies, so convergence is guaranteed in a finite number of iterations



1. Initialization
   $V(s) \in \Re$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Repeat
   $\Delta \leftarrow 0$
   For each $s \in \mathcal{S}$:
   $v \leftarrow V(s)$
   $V(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' \mid s, \pi(s)) V(s')$
   $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $b \leftarrow \pi(s)$
   $\pi(s) \leftarrow \arg\max_{a \in A} \left\{ R(s, a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right\}$
   If $b \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop; else go to 2

# Policy iteration discussion

- There is a finite number of stationary policies:
  - A total of $|S|$ states and $|A|$ actions possible in each state $\rightarrow$ A total of $|A|^{|S|}$ unique stationary policies

- Policy iteration converges in a finite number of steps

- Policy evaluation can be done iteratively (as shown before) or by solving a system of linear equations

- When the state space is large, this is expensive

# Relationship between policy iteration and value iteration

Recall value iteration...

Initialize $V$ arbitrarily, e.g., $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat
    $\Delta \leftarrow 0$
    For each $s \in \mathcal{S}$:
        $v \leftarrow V(s)$
        $V(s) \leftarrow \max_{a \in A} \left\{ R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right\}$
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, $\pi$, such that
    $\pi(s) \leftarrow \arg\max_{a \in A} \left\{ R(s,a) + \gamma \sum_{s' \in S} p(s' \mid s, a) V(s') \right\}$

Value iteration truncates policy iteration by combining one sweep of policy evaluation and one of policy improvement in each of its sweeps.

# Relationship between policy iteration and value iteration

- Each iteration of value iteration is relatively cheap compared to iterations of policy iteration because policy iteration requires solving a system of $|S|$ linear equations in each iteration.

- However, the trade off is that policy iteration requires less iterations to converge.

- Implementations in MATLAB
  - **Value iteration:** "slow" because MATLAB is bad at iterative computations
  - **Policy iteration:** "fast" because MATLAB is very good at solving systems of linear equations

# Complete knowledge

- In problems of *complete knowledge*, the agent has a complete and accurate model of the environment's dynamics.

- If the environment is an MDP, then such a model consists of the one-step *transition probabilities* and *expected rewards* for all states and their allowable actions.

- In problems of *incomplete knowledge*, a complete and perfect model of the environment is not available.

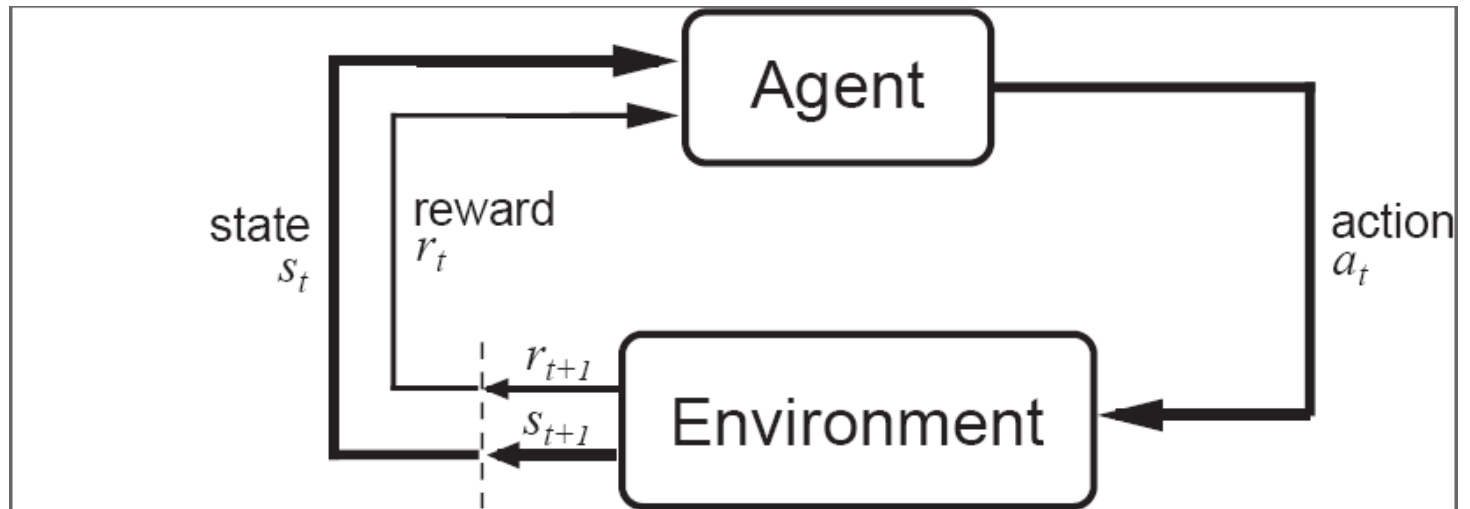# What happens if the environment is unknown?

- Model-based
  - Learn reward and transition probabilities
  - Then compute optimal value function and corresponding policy
  - Example: RTDP (real-time dynamic programming)

- Model-free
  - Learn value function or action value function directly

# Why is Online Learning Important?

- If the state transition probabilities are known, finding the optimal policy becomes a straightforward computational problem, however…

- If the transition probabilities and/or rewards are unknown, then this is a *reinforcement learning* (RL) problem.

- Reinforcement-learning techniques learn from experience, no knowledge of the model is required.

# Agent-environment interaction

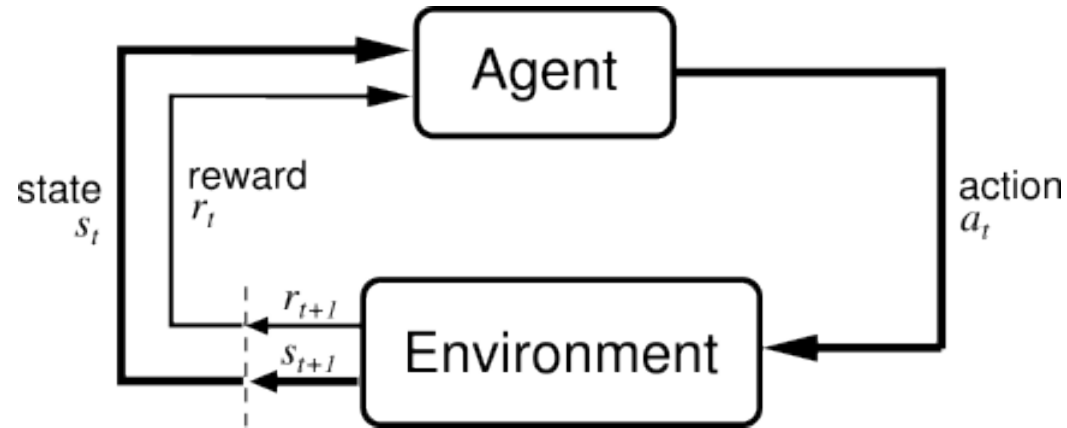By "state" we mean whatever information is available to the agent



Everything inside the agent is completely known and controllable by the agent; everything outside is incompletely controllable but may or may not be completely known.

In our setting, who are the agent and the environment? What are the states, actions, rewards?

# Typical Agent

- In reinforcement learning (RL), the agent observes a state and takes an action.

- Afterward, the agent receives a reward.



Goal: Optimize Cumulative Reward

- Rewards are calculated in the environment

- Used to teach the agent how to reach a goal state

- Must signal what we ultimately want achieved, not necessarily subgoals

- May be *discounted* over time

- In general, seek to maximize the *expected return*

*Examples of reinforcement learning techniques: Q-learning, Sarsa, actor critic etc.*

# Estimation of Action Values

- State values $V^{\pi}(s)$ are not enough without a model
- We need action values as well:

   $Q^{\pi}(s, a)$ → expected return when starting in state $a$, taking action $a$, and thereafter following policy $\pi$

- Exploration vs. Exploitation
   – The agent believes that a particular action has a high value; should it choose that action all the time, or should it choose another one that it has less information about, but seems to be worse?

# Another Value Function: State-action-value Function

- $Q^{\pi}$ (also called Q-function) is the expected return (value) starting from state $s$, taking action $a$, and thereafter following policy $\pi$
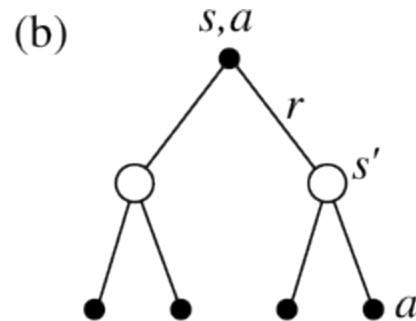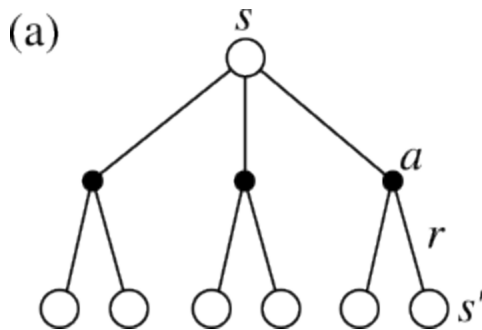
State-action-value function for policy $\pi$

$$Q^{\pi}(s,a) = E_{\pi}\left[\sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s, a_0 = a\right]$$

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

$$V^*(s) = \max_{a \in \mathcal{A}} Q^*(s,a) \qquad \pi^*(s) = \arg\max_{a \in \mathcal{A}} Q^*(s,a)$$



Backup diagrams for (a) $V^{\pi}$ and (b) $Q^{\pi}$

# Q-learning

- Q-learning algorithm does not specify what the agent should actually do.

- The agent learns a *Q*-function that can be used to determine an optimal action.

- There are two things that are useful for the agent to do:

  - **exploit** the knowledge that it has found for the current state *s* by chosing one of the actions *a* that maximizes *Q[s,a]*.

  - **explore** in order to build a better estimate of the optimal *Q*-function. That is, it should select a different action from the one that it currently thinks is best.

# Q-learning

- At each step (denote by s the current state), choose the action a which maximizes the function Q(s, a)

- Q is the estimated utility function – it tells us how good an action is given a certain state (include the estimated future utility)

- Q(s, a) = immediate reward for making an action + best utility (Q) for the resulting state

- Q learning:

$$Q(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha_t}_{\text{learning rate}} \cdot \left( \overbrace{\underbrace{r_{t+1}}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}}}^{\text{learned value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}} \right)$$

where $r_{t+1}$ is the reward observed after performing $a_t$ in $s_t$, and where $\alpha_t$ is the learning rate ($0 < \alpha_t \leq 1$).

# Q-learning discussion

- $Q$-learning is guaranteed to converge to the optimal $Q$-values if all $Q(s, a)$ values are updated infinitely often Watkins and Dayan (1992).

- In order to make sure all actions will eventually be tried in all states exploration is necessary.

- A common exploration method is to execute a random action with small probability $\epsilon$, which is known as $\epsilon$-greedy exploration Sutton and Barto (1998).

# Other reinforcement learning techniques

- Sarsa, Actor critic, etc.
- However, this goes beyond the topic of this class

# Problems with Generic Learning Algorithms

- Curse of Dimensionality (often multi-dimensional state space)

- Slow convergence (not useful for real-time applications)

- Solutions to these problems exist