

MARKOV DECISION PROCESSES (MDP) AND REINFORCEMENT LEARNING (RL)

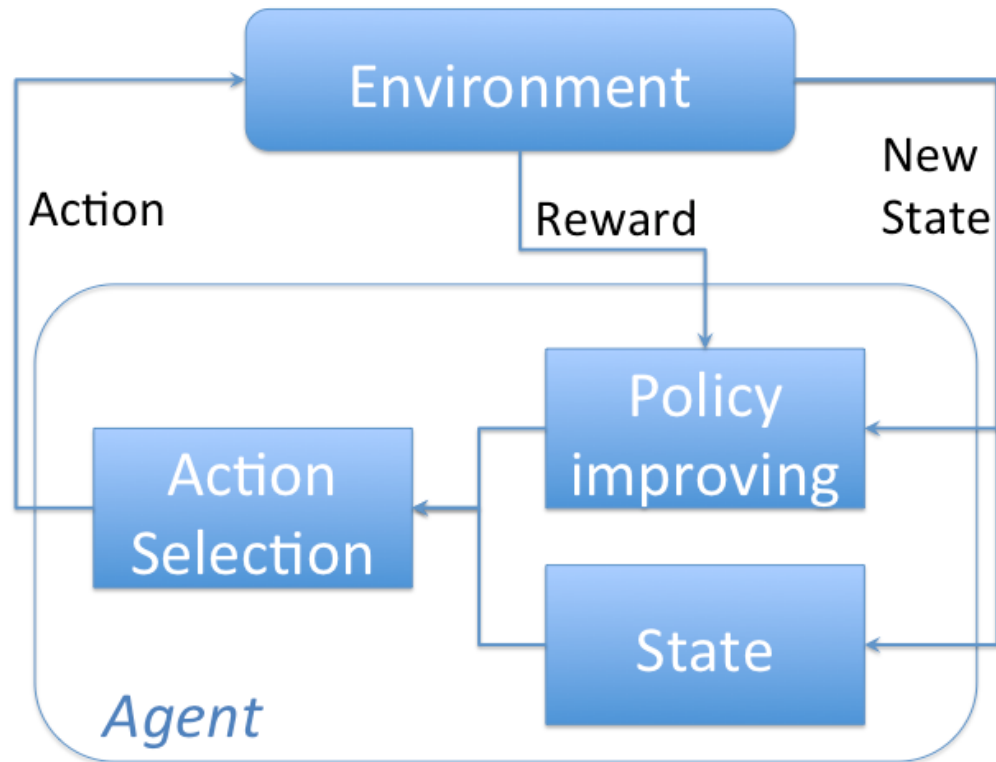
Historical background

2

- Original motivation: animal learning
- Early emphasis: neural net implementations and heuristic properties
- Now appreciated that it has close ties with
 - ▣ operations research
 - ▣ optimal control theory
 - ▣ dynamic programming
 - ▣ AI state-space search
- Best formalized as a set of techniques to handle: Markov Decision Processes (MDPs) or Partially Observable Markov Decision Processes (POMDPs)

Reinforcement learning task

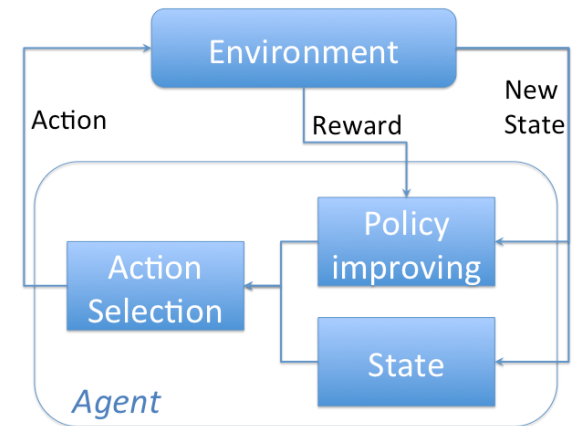
3



Reinforcement learning task

4

$$s(0) \xrightarrow[r(1)]{a(1)} s(1) \xrightarrow[r(2)]{a(2)} s(2) \xrightarrow[r(3)]{a(3)} s(3) \dots$$



- **Goal:** learn to choose actions that maximize the **cumulative** reward

$$r(0) + \gamma r(1) + \gamma r(2) + \gamma r(3) + \dots$$

where $0 \leq \gamma < 1$

Foresighted Optimization

5

- The key feature of this framework is that actions affect immediate and future system performance
 - ▣ We optimize the system “on the long run”
- In scenarios in which actions/decisions affect immediate and future performance, myopic heuristic solutions are suboptimal because they ignore the expected future utility
- Dramatic improvements can be achieved using *long term* optimization

Stochastic Process

6

- Quick definition: a random process
- Often viewed as a collection of indexed random variables $\{X_t: t \in T\}$
 - $X_t \in S$ is the system state
 - t is the time
- Useful to characterize “environment” dynamics
 - Set of states with probability law governing evolution of states over time
- We will focus on *discrete-time stochastic chains*
 - Time is discrete
 - State set S is discrete

Stochastic Process Dynamics

- The process *dynamics* can be defined using the transition probabilities
- They specify the stochastic evolution of the process through its states
- For a discrete time process, transition probabilities can be defined as follows

$$P(X_{t+1} = x_{t+1} | X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_0 = x_0)$$

Markov Property

- The term **Markov property** refers to the memoryless property of a stochastic process:
- For a discrete time process, the Markov property is defined as:

$$\begin{aligned} P(X_{t+1} = x_{t+1} | X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_0 = x_0) \\ = \\ P(X_{t+1} = x_{t+1} | X_t = x_t) \end{aligned}$$

- **Definition:** a stochastic process that satisfies the Markov property is called Markov process
- If the state space is discrete, we refer to these processes as **Markov Chains**

Markov Property

9

- Why useful?
 - ▣ Simple model of temporal correlation of environment dynamics
 - ▣ Current state contains all information needed to predict distribution of future state(s)
- Does it hold in the real world?
- It's an ideal
 - ▣ Will allow us to prove properties of algorithms
 - ▣ Algorithms mostly still work

Markov Chain

10

- Let $\{X_t: t = 1, 2, 3, \dots\}$ be a Markov chain
- Dynamics is defined using the transition probability function

$$P(X_{t+1} = s' | X_t = s) \\ t \geq 0, \quad s', s \in S$$

- Under some assumptions (see previous lesson), a Markov chain has a *stationary* transition probability function

$$\pi_j = \lim_{k \rightarrow \infty} \pi_j(k)$$

Markov Decision Processes

11

- Finite set of states S
- Finite set of actions $A(s), s \in S$
- Immediate reward function
$$R: S \times A \rightarrow R$$
- Transition (next-state) function
$$T: S \times A \rightarrow S$$
- More generally, R and T are treated as stochastic
 - ▣ We'll stick to the above notation for simplicity
 - ▣ In general case, treat the immediate rewards and next states as random variables, take expectations, etc.
- We will focus on *discrete time* MDPs

Markov Decision Processes

12

- Markov Property for MDPs

$$P(s' | s, a)$$
$$s', s \in S, \quad a \in A(s)$$

- Next state is a function of current state and the action taken!

Markov Decision Processes

13

- If no rewards and only one action, this is just a *Markov chain*
- Sometimes also called a **Controlled Markov Chain**
- Overall objective is to determine a policy

$$\pi : S \rightarrow A$$

such that some measure of cumulative reward is optimized

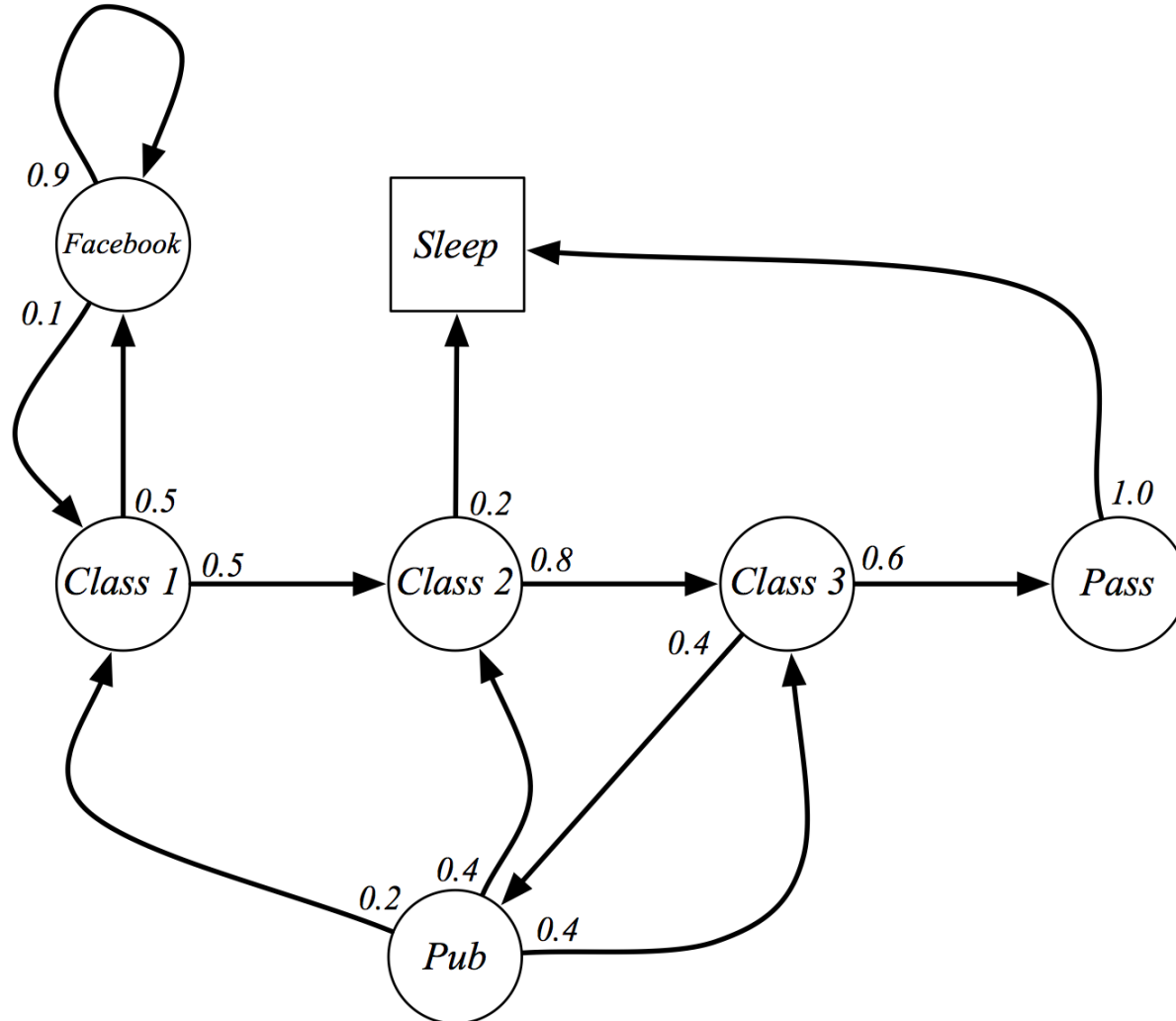
- ▣ E.g, $r(0) + \gamma r(1) + \gamma r(2) + \gamma r(3) + \dots$

What's a policy?

14

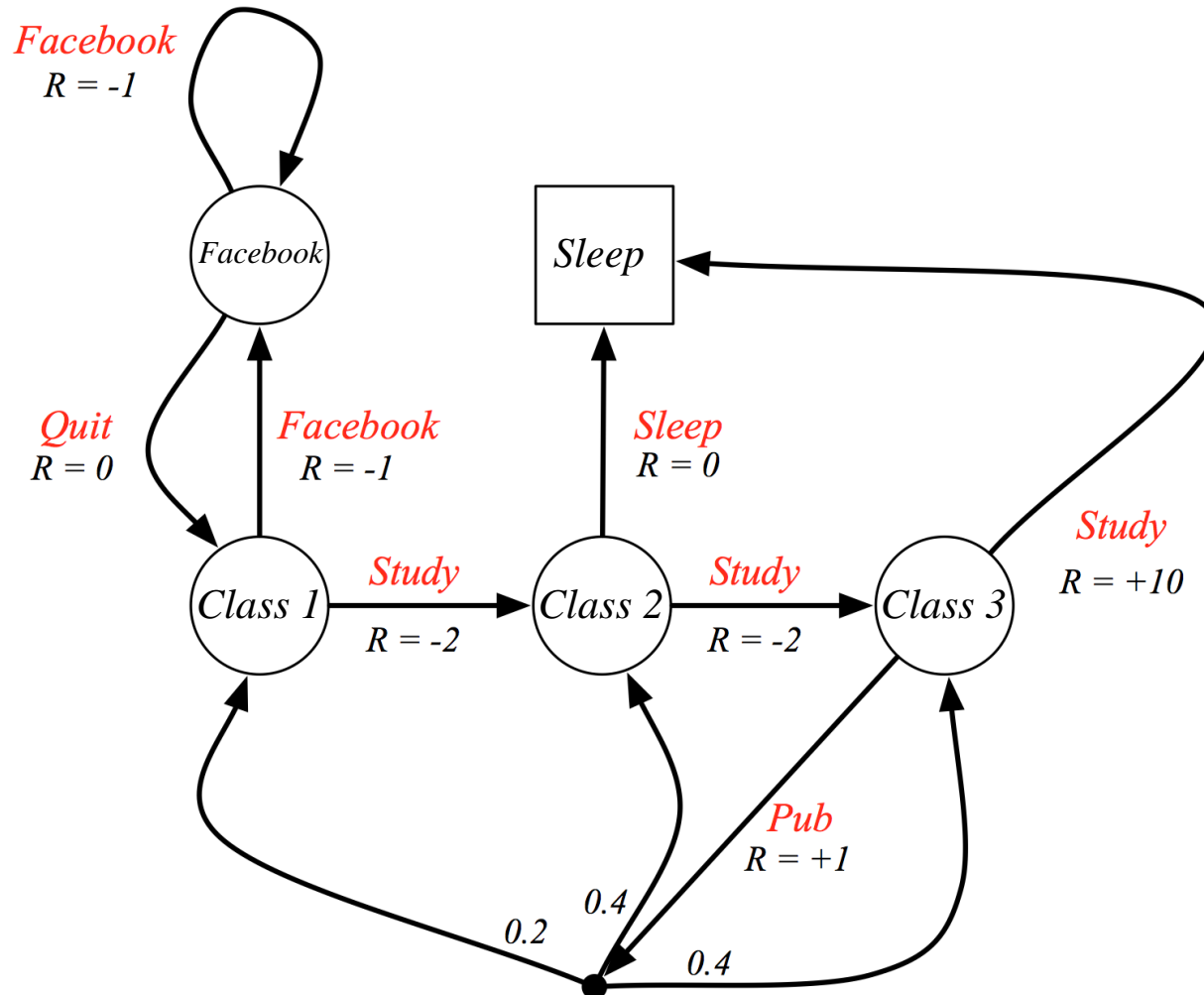
If agent is in this state	Then a good action is
s_1	a_3
s_2	a_7
s_3	a_1
s_4	a_3
...	...

Student Markov Chain



Student Markov Decision Process

16

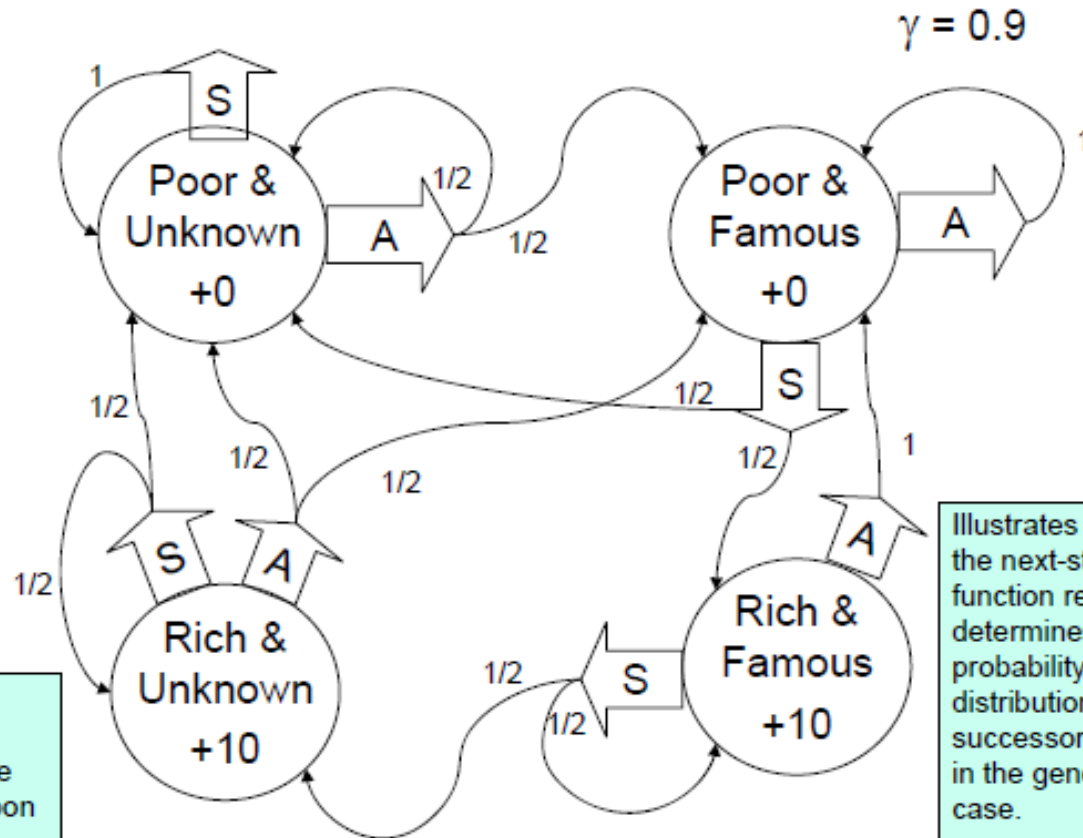


A Markov Decision Process

17

You run a startup company. In every state you must choose between Saving money or Advertising.

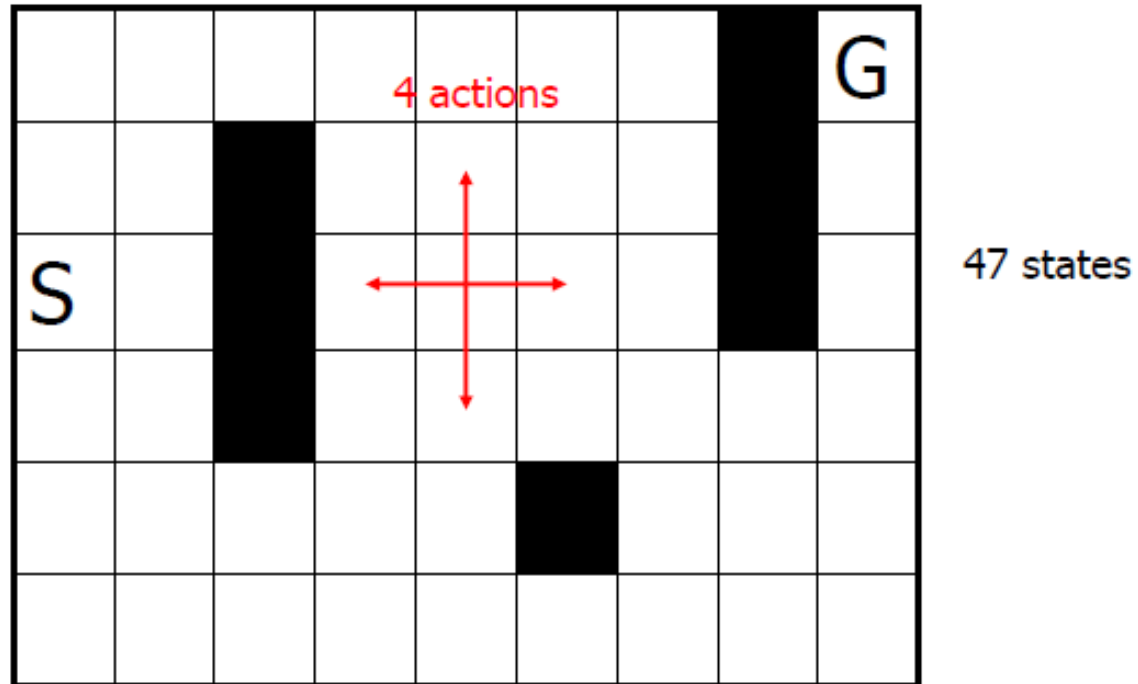
Here the reward shown inside any state represents the reward received upon entering that state.



Illustrates that the next-state function really determines a probability distribution over successor states in the general case.

Another MDP

18



Reward = -1 at every step

$\gamma = 1$

G is an absorbing state, terminating any single trial, with a reward of 100

Effect of actions is deterministic

Applications of MDPs

19

Many important problems are MDPs....

- ... Robot path planning
- ... Travel route planning
- ... Elevator scheduling
- ... Autonomous aircraft navigation
- ... Manufacturing processes
- ... Network switching & routing

And many of these have been successfully handled using RL methods

Brief summary of concepts

20

- 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

Value Functions

21

□ It turns out that

□ RL theory

□ MDP theory

□ AI game-tree search

all agree on the idea that evaluating states is a useful thing to do.

□ A **(state) value function** V is any function mapping states to real numbers:

$$V : S \rightarrow \mathbb{R}$$

What's a value function?

22

If agent starts in this state	Return when following given policy should be
s_1	13
s_2	-1
s_3	22.6
s_4	6
...	...

A special value function: the return

23

- For any policy π , define the **return** to be the function $V^\pi: S \rightarrow \mathbb{R}$ assigning to each state the quantity

$$V^\pi(s) = \sum_{t=0}^{\infty} \gamma^t r(t)$$

Reminder: Use expected values in the stochastic case.

where

- $s(0) = s$
- each action $a(t)$ is chosen according to policy π
- each subsequent $s(t + 1)$ arises from the transition function T
- each immediate reward $r(t)$ is determined by the immediate reward function R
- γ is a given discount factor in $[0, 1]$

Why the discount factor γ ?

24

- Models idea that future rewards are not worth as much as immediate rewards
 - Used in economic models
 - Uncertainty about the future
- Also models situations where there is a nonzero fixed probability $1 - \gamma$ of termination at any time
- Tradeoff between myopic ($\gamma = 0$) vs foresighted optimization (γ close to 1)
- ...and makes the math work out nicely
 - with bounded rewards, sum guaranteed to be finite even in infinite-horizon case

Technical remarks

25

- If the next state and/or immediate reward functions are stochastic, then the $r(t)$ values are random variables and the return is defined as the expectation of this sum
- If the MDP has absorbing states, the sum may actually be finite
 - ▣ In that case $\gamma = 1$ is allowed, i.e., no discount
 - ▣ We stick with this infinite sum notation for the sake of generality
 - ▣ The formulation we use is called infinite-horizon

Optimal Policies

26

- Objective: Find a policy π^* such that

$$V^{\pi^*}(s) \geq V^{\pi}(s)$$

for any policy π and any state s

- Such a policy is called an **optimal** policy

- We define:

$$V^* = V^{\pi^*}$$

Interesting fact

27

- For every MDP such that S is discrete and $A(s)$ is finite there exists an optimal policy
 - ▣ This theorem can be easily extended to the case in which $A(s)$ is a *compact* set
- It's a policy such that for every possible start state there is no better option than to follow the policy

Finding an Optimal Policy

28

- Idea One:

 - Run through all possible policies.

 - Select the best.

- What's the problem ??

Finding an Optimal Policy

29

- Dynamic Programming approach:
 - Determine the optimal value function for each state
 - Select actions according to this optimal value function V^*

- How do we compute V^* ?
 - Magic words: **Bellman equation(s)**

Simple derivation of the Bellman equation

30

- Given the state transition $s \rightarrow s'$

$$\begin{aligned} V^\pi(s) &= \sum_{t=0}^{\infty} \gamma^t r(t) \\ &= r(0) + \gamma \sum_{t=0}^{\infty} \gamma^t r(t+1) \\ &= r(0) + \gamma V^\pi(s') \end{aligned}$$

Bellman equations

31

- For any state s and policy π

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

- For any state s , the optimal value function is

$$V^*(s) = \max_a \{R(s, a) + \gamma V^*(T(s, a))\}$$

- **Recurrence relations**

- ▣ Can be used to compute the return from a given policy or to compute the optimal return via value iteration

Bellman equations: general form

32

- For completeness, here are the Bellman equations for stochastic and discrete time MDPs:

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} P_{ss'}(\pi(s)) V^\pi(s')$$
$$V^*(s) = \max_a \{R(s, a) + \gamma \sum_{s'} P_{ss'}(a) V^*(s')\}$$

where $R(s, a)$ now represents $E(R | s, a)$ and

$P_{ss'}(a)$ = probability that the next state is s'
given that action a is taken in state s

From values to policies

34

- Given the optimal value function V^* it follows from Bellman equation that the optimal policy can be computed as:

$$\pi(s) = \underset{a}{\operatorname{argmax}}\{R(s, a) + \gamma V^*(s')\}$$

- An optimal policy is said to be greedy for V^*
- If π is not optimal then a greedy policy for V^π will yield a larger return than π
 - ▣ Not hard to prove
 - ▣ Basis for another DP approach to finding optimal policies: policy iteration

Finding an optimal policy

44

Value Iteration Method

Choose any initial state value function V_0

Repeat for all $n \geq 0$

For all s

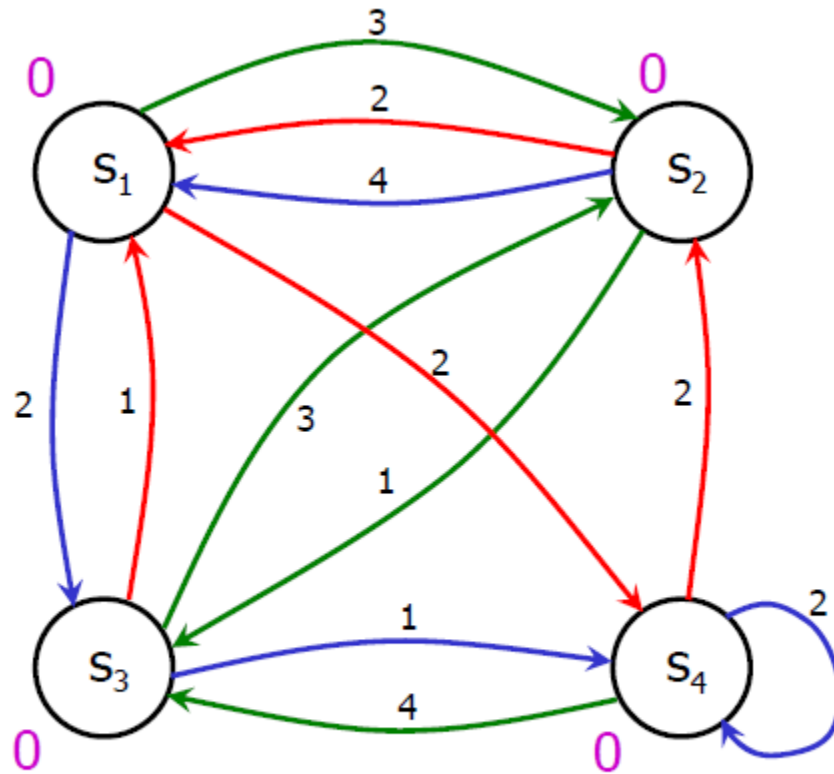
$$V_{n+1}(s) \leftarrow \max_a \{R(s, a) + \gamma V_n(T(s, a))\}$$

Until convergence

- This converges to V^* and any greedy policy with respect to it will be an optimal policy

Value Iteration

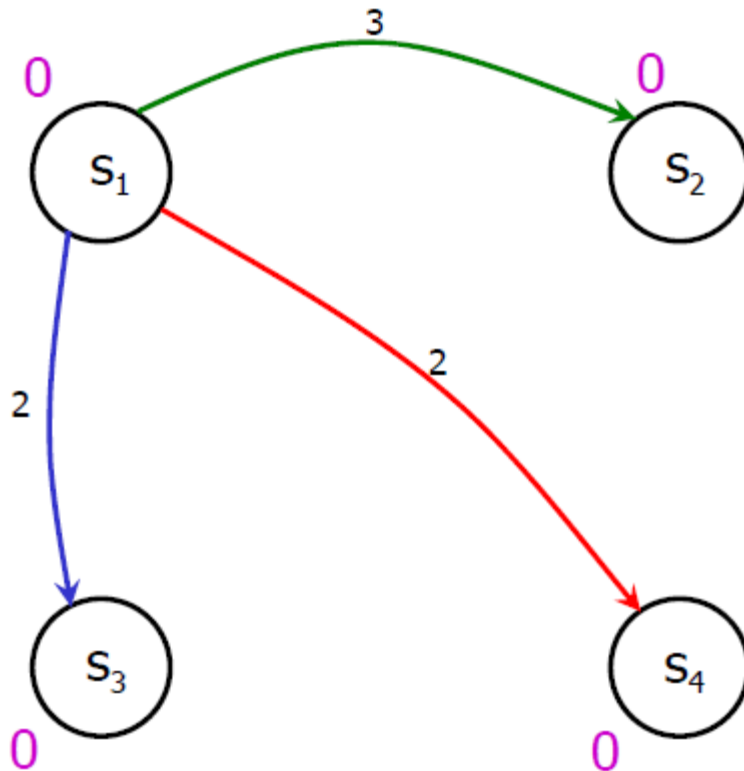
45



Arbitrary initial value function V_0

Value Iteration

46



Arbitrary initial value function V_0

Computing a new value for s_1 using 1-step lookahead with previous values:

For action a_1 lookahead value is $2 + (.9)(0) = 2$

For action a_2 lookahead value is $3 + (.9)(0) = 3$

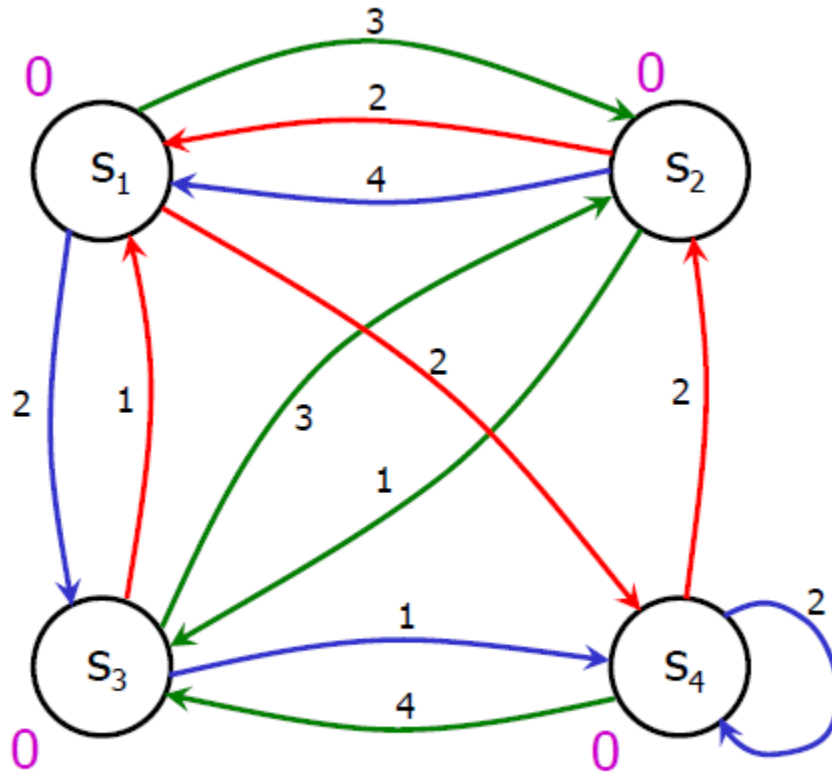
For action a_3 lookahead value is $2 + (.9)(0) = 2$

a_1	a_2	a_3
2	3	2

$$V_1(s_1) = \max\{2, 3, 2\} = 3$$

Value Iteration

47

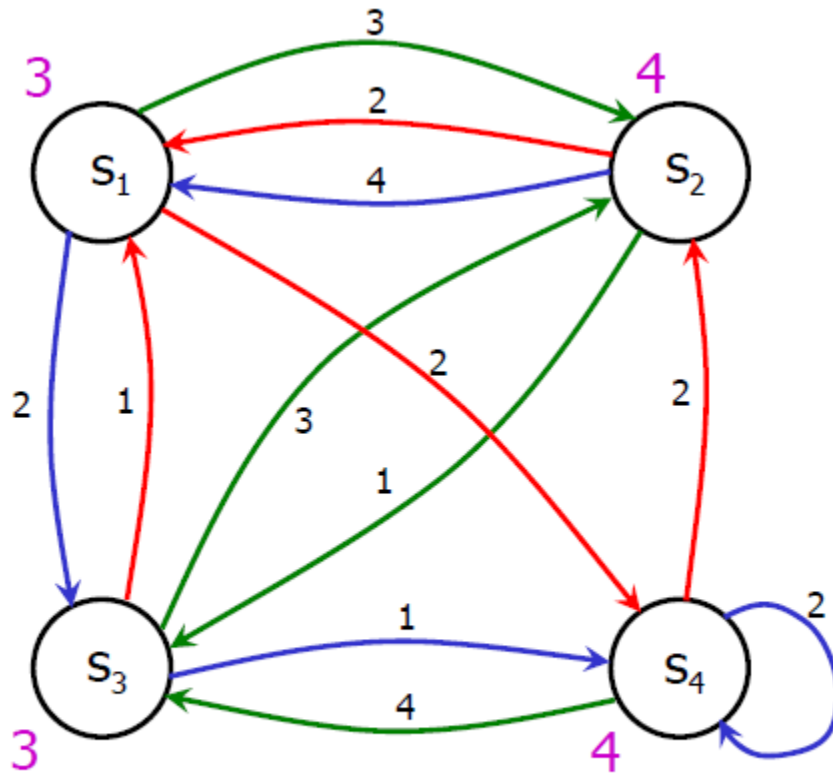


Arbitrary initial value function V_0

Lookahead value along action				
	a_1	a_2	a_3	max
s_1	2	3	2	3
s_2	2	1	4	4
s_3	1	3	1	3
s_4	2	4	2	4

Value Iteration

48



Updated
approximation
to V^* :

$$V_1(s_1) = 3$$

$$V_1(s_2) = 4$$

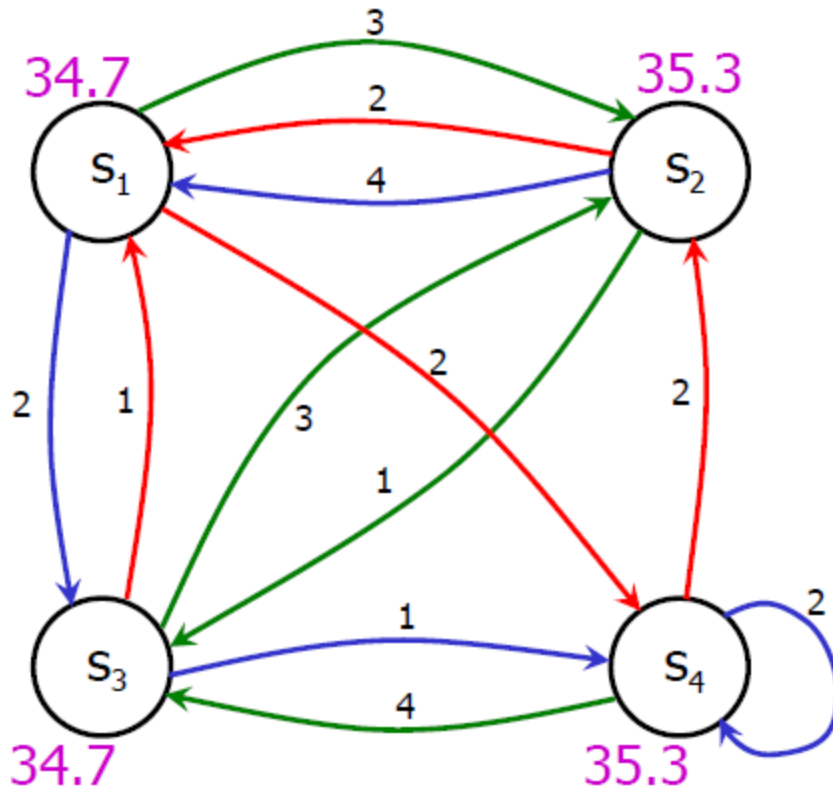
$$V_1(s_3) = 3$$

$$V_1(s_4) = 4$$

New value function V_1 after one step of value iteration

Value Iteration

49

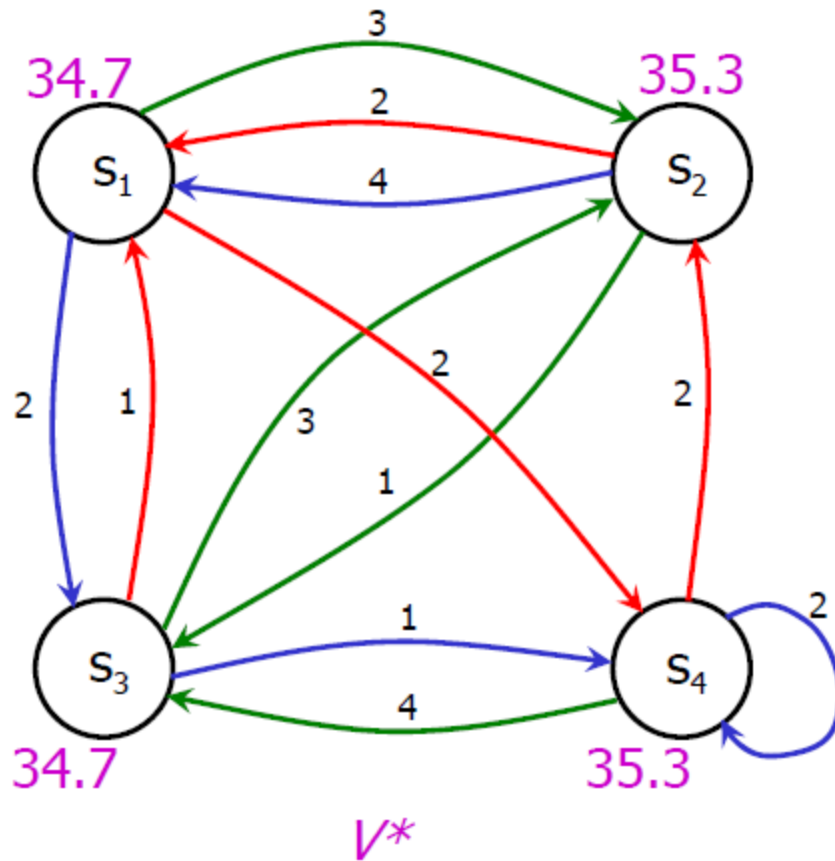


	s_1	s_2	s_3	s_4
V_0	0	0	0	0
V_1	3	4	3	4
V_2	6.6	6.7	6.6	6.7
V_3	9.0	9.9	9.0	9.9
V_4	11.9	12.1	11.9	12.1
V_5	13.9	14.8	13.9	14.8
...				
V^*	34.7	35.3	34.7	35.3

Keep doing this until it converges to V^*

Value Iteration

50

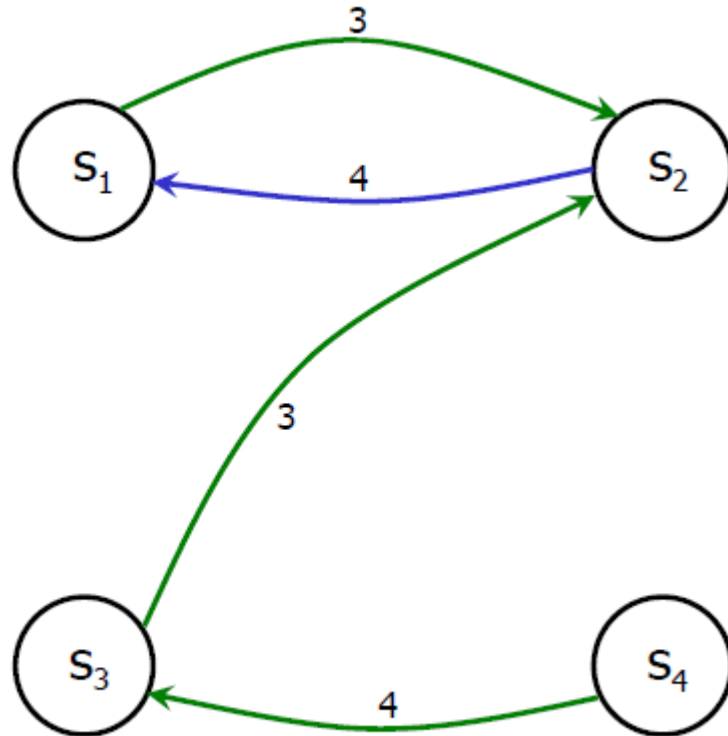


Determining a greedy policy for V^*

Lookahead value along action				
	a_1	a_2	a_3	best
s_1	33.8	34.8	33.2	a_2
s_2	33.2	32.2	35.2	a_3
s_3	32.2	34.8	32.8	a_2
s_4	33.8	35.2	33.8	a_2

Value Iteration

51



Optimal policy

Value iteration – Full Version

52

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

Repeat

$$\Delta \leftarrow 0$$

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

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$$

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

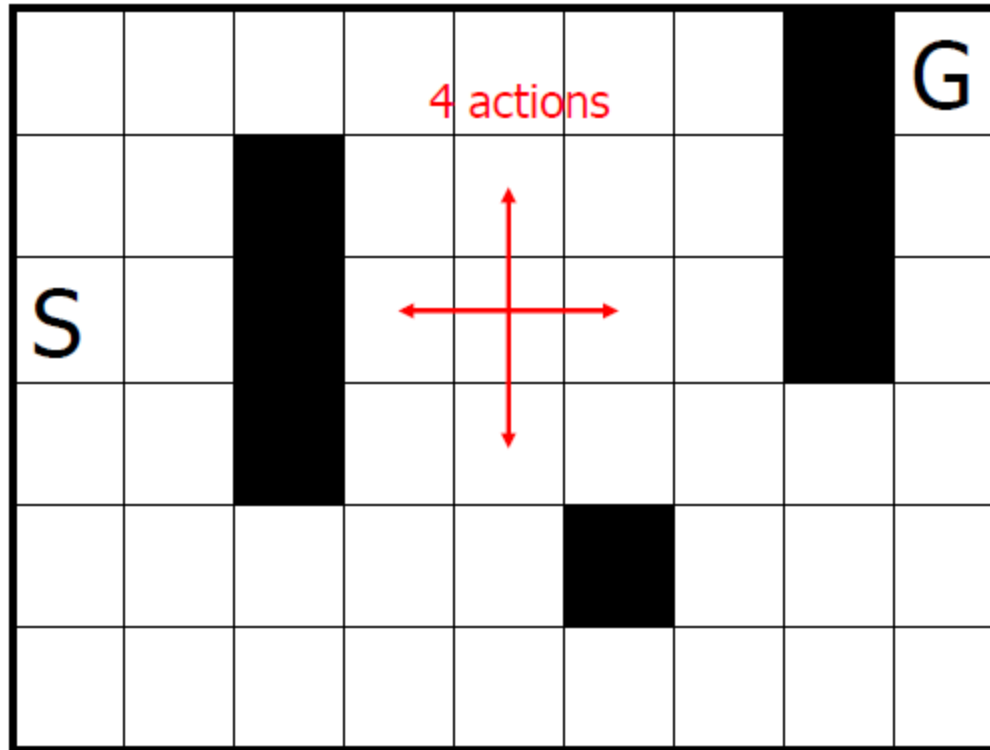
until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, π , such that

$$\pi(s) = \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$$

Maze Task

53



Reward = -1 at every step

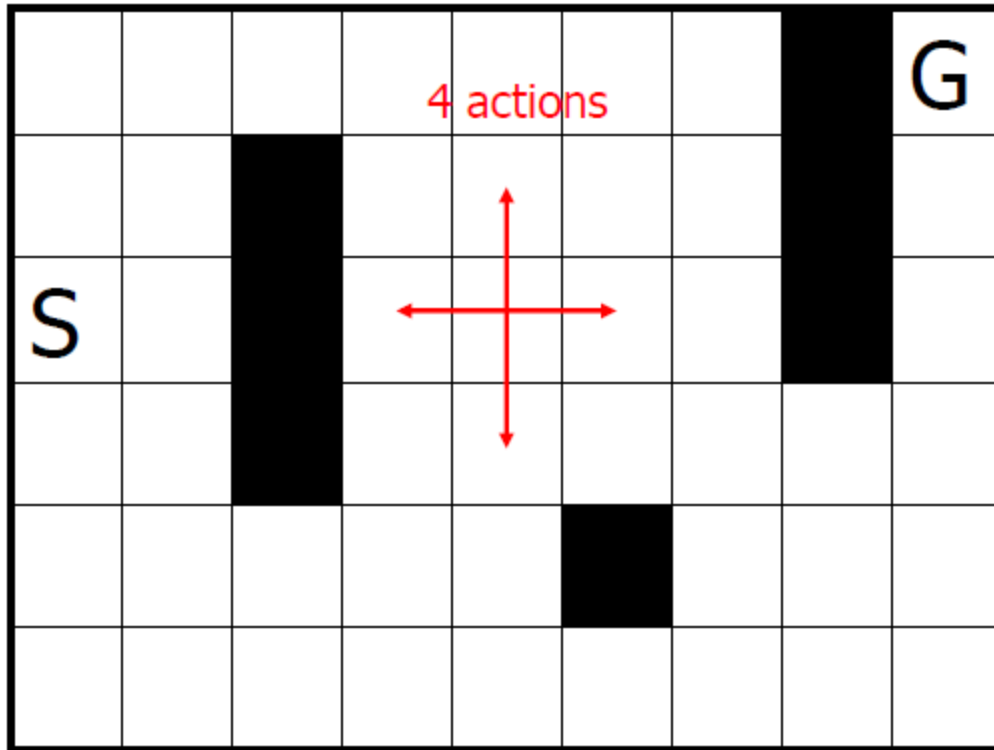
$\gamma = 1$

G is an absorbing state, terminating any single trial, with a reward of 100

Effect of actions is deterministic

Maze Task

54



How would you model the MDP?

Reward = -1 at every step

$\gamma = 1$

G is an absorbing state, terminating any single trial, with a reward of 100

Effect of actions is deterministic

Maze Task – MDP model

55

- State is a couple: $s = (x, y)$, $x, y \in \{1, \dots, 9\}$ defining the robot position
- Actions: $A(s) = \{\text{up, down, left, right}\}$
(except for those states near the black squares)
- Reward Function: $R(s) = \begin{cases} -1, & \forall s \neq G \\ 100, & \text{if } s = G \end{cases}$
- Transition function: $s' = \begin{cases} (x + 1, y) & \text{if } a = \text{right} \\ (x - 1, y) & \text{if } a = \text{left} \\ (x, y + 1) & \text{if } a = \text{up} \\ (x, y - 1) & \text{if } a = \text{down} \end{cases}$

Maze Task – Value function

56

	86	87	88	89	90	91	92		100	G
	85	86		90	91	92	93		99	
S	86	87		91	92	93	94		98	
	87	88		92	93	94	95	96	97	
	88	89	90	91	92		94	95	96	
	87	88	89	90	91	92	93	94	95	

V^*

What's an optimal path from S to G?

Maze Task – Optimal Path

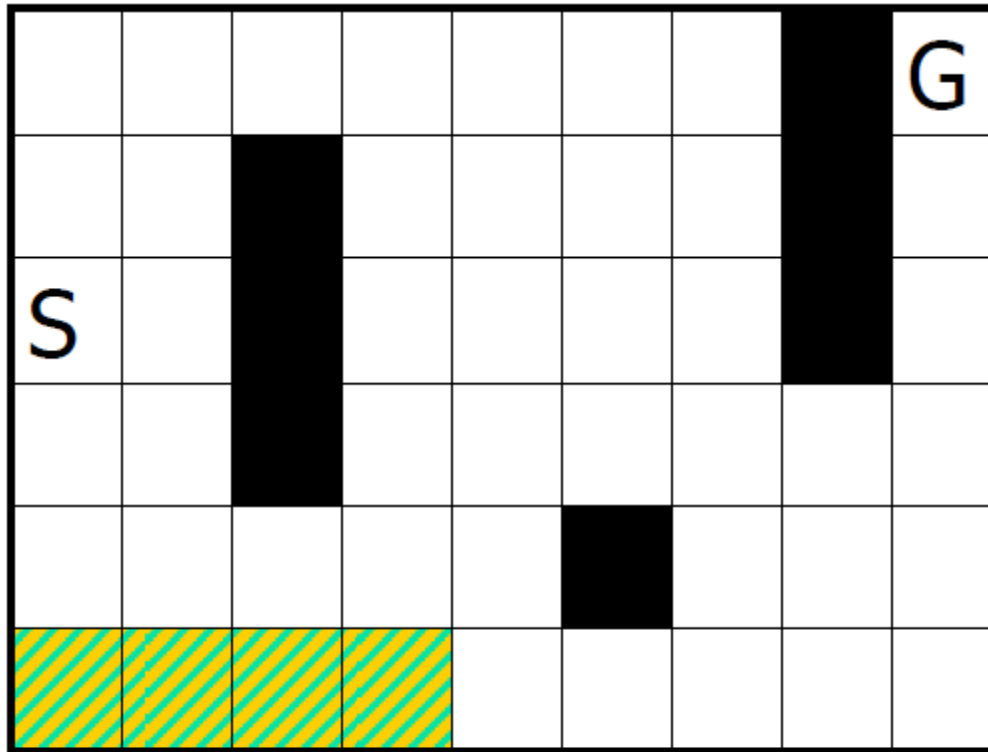
57

	86	87	88	89	90	91	92		100	G
	85	86		90	91	92	93		99	
S	86	87		91	92	93	94		98	
	87	88		92	93	94	95	96	97	
	88	89	90	91	92		94	95	96	
	87	88	89	90	91	92	93	94	95	

V^*

Another Maze Task

58



Now what's an optimal path from S to G?

With:

$P=0.1$ to the right

$P=0.1$ to the left

Everything else same as before, except:

With some nonzero probability, a small wind gust might displace the agent one cell to the right or left of its intended direction of travel on any step

Entering any of the 4 patterned cells at the southwest corner yields a reward of -100

Another Maze Task – MDP model

59

- Reward function: same as before, except that

$$R(s) = -100, \quad \forall s \in S: x \leq 4, y = 1$$

- Transition function:

$$s = (x, y), a = \text{up}$$

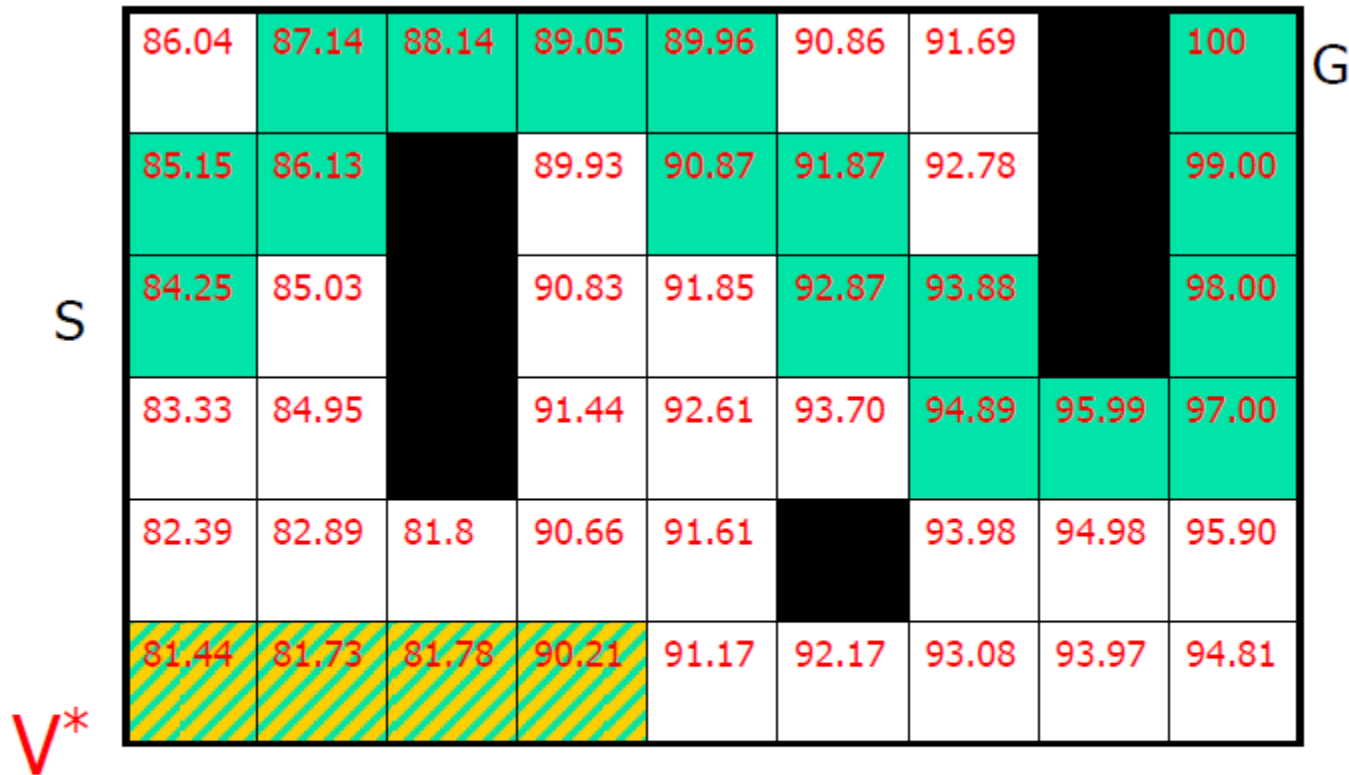
$$P(s' = (x, y + 1) \mid s, a) = 0.8$$

$$P(s' = (x + 1, y) \mid s, a) = 0.1$$

$$P(s' = (x - 1, y) \mid s, a) = 0.1$$

Another Maze Task

60



With probability 0.2, a small wind gust might displace the agent one cell to the right or left of its intended direction of travel on any step

Entering any of the 4 patterned cells at the southwest corner yields a reward of -100

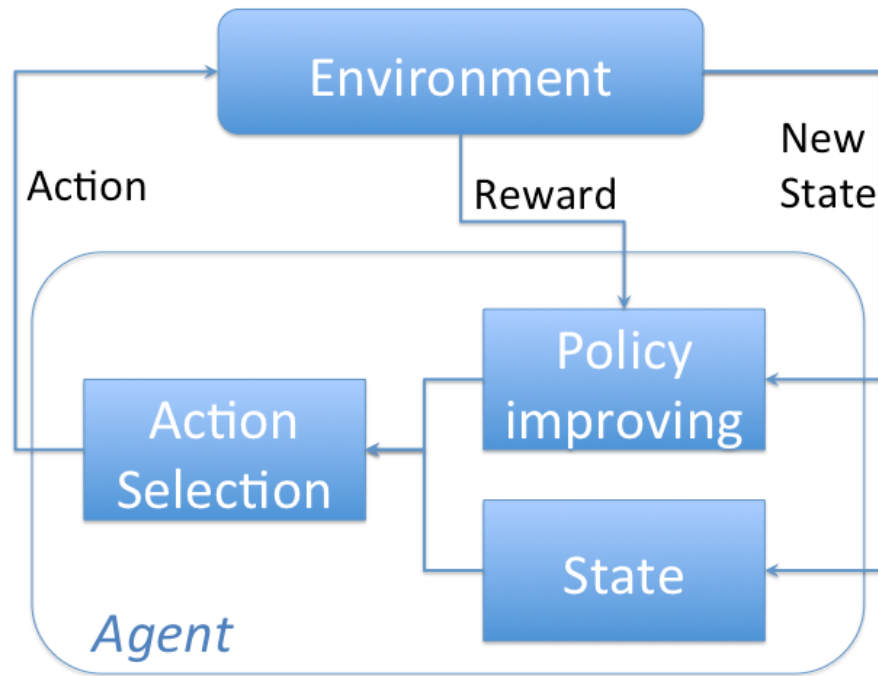
Why is On-line Learning Important?

61

- So far, we assumed all system parameters to be known
- Finding the optimal policy becomes a straightforward computational problem
 - ▣ E.g., value iteration, but even policy iteration, linear programming, ecc...
- What if rewards/transitions probabilities are unknown? Can we compute the optimal policy?
- We have to deal with a *Reinforcement Learning problem!*

Agent-Environment Interaction

62



- 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

Agent Knowledge

63

- A reinforcement learning problem can be posed in a variety of different ways depending on assumptions about the level of knowledge initially available to the agent
- 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

Q-Values

66

- For any policy π , define $Q^\pi: S \times A \rightarrow R$ by

$$Q^\pi(s, a) = \sum_{t=0}^{\infty} \gamma^t r(t)$$

Again, the correct expression for a general MDP should use expected values here

- $s(0) = s$ is the initial state,
 - $a(0) = a$ is the action taken,
 - and all subsequent states, actions, and rewards arise following policy π
-
- Just like V^π except that action a is taken at the very first step and only after this, policy π is followed
 - Bellman equations can be rewritten in terms of Q-values

What are Q-values?

67

If agent is in this state	And starts with this action and then follows the policy	Return should be
s_1	a_1	-5
s_1	a_2	3
s_2	a_1	17.1
s_2	a_2	10
...

Q-Values

68

- Relationship between Value function and Q-function (given the state transition $s \rightarrow s'$)

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

versus

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

Q-Values

69

- Relationship between Value function and Q-function

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

versus

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

Q-Values

70

- Define $Q^* = Q^{\pi^*}$, where π^* is an optimal policy

$$Q^*(s, a) = R(s, a) + \gamma V^*(s')$$

- Since:

$$V^*(s) = \max_a \{R(s, a) + \gamma V^*(s')\}$$

- Then:

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

- And:

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

Q-Values

71

- The optimal policy π^* is greedy for Q^* , that is

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

[it follows from $V^*(s) = \max_a Q^*(s, a)$]

Q-learning Algorithm

72

- Q is the estimated utility function
 - ▣ It tells us how good an action is, given a certain state
 - ▣ It includes immediate reward for making an action + best utility (Q) for the resulting state (future utility)
 - ▣ It allows to compute the optimal policy
- Q-learning is based on an online estimation of the Q function

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha [r(s, a) + \gamma \max_{a'} Q(s', a')]$$

Q-learning Algorithm

73

Initialize $Q(s, a)$ arbitrarily

Repeat (for each decision epoch)

Initialize s

Repeat (for each step of episode)

Choose a from s using a policy derived from Q

Take action a , observe $r(s, a)$,

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha [r(s, a) + \gamma \max_{a'} Q(s', a')]$$

$$s \leftarrow s'$$

until s is terminal

Exploitation and exploration

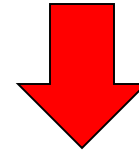
- Q-learning algorithm does not specify what the agent should actually do. The agent learns a Q-function that can be used to determine an optimal action. There are two things that are useful for the agent to do:
 - **exploit** the knowledge that it has found for the current state s by doing one of the 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.

Exploitation and exploration

75

Choose a from s using a policy derived from Q

- Simple Approach: ϵ -greedy policy
 - ϵ small number, e.g., 0.1



Generate a random number p

if $p \leq \epsilon$

Choose an action at random \rightarrow **explore**

else

Choose the greedy action $a^* = \arg \max_a Q(s, a) \rightarrow$ **exploit**

end

Q-learning discussion

76

- 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)
- It follows that exploration is necessary
 - ▣ A common approach is the ϵ -greedy strategy
- Q-learning can be very slow to converge to the optimal policy, especially if the state space is large
- One of the biggest challenges in the RL field is to *speed up* the learning process

Learning or planning?

77

- Classical DP emphasis for optimal control
 - ▣ Dynamics and reward structure known
 - ▣ Off-line computation
- Traditional RL emphasis
 - ▣ Dynamics and/or reward structure initially unknown
 - ▣ On-line learning
- Computation of an optimal policy off-line with known dynamics and reward structure can be regarded as planning

More information

78

- You can find more information about MDPs and learning here:

<https://webdocs.cs.ualberta.ca/~sutton/book/the-book.html>