

# The Learning Task

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in The Learning Task

1. The Reinforcement Learning Task
2. Markov Decision Process (MDP)
3. Discounted Cumulative Reward
4. Example of Simple grid-world

# The Learning Task

- Formulate the problem of learning sequential control strategies more precisely
- Many ways to do so
  - Agents actions are deterministic or non-deterministic
  - Agent can predict next state due to an action or it cannot
  - Training by an expert may show optimal action sequences
- A general formulation is based on Markov Decision Processes

# Markov Decision Process (MDP)

- Agent can perceive a set  $S$  of discrete states in its environment
- Has a set  $A$  of actions that it can perform
- At each discrete time step  $t$  the agent senses the current state  $s_t$ , chooses a current action  $a_t$  and performs it
- Environment responds by giving it a reward  $r_t = r(s_t, a_t)$  and by producing state  $s_{t+1} = \delta(s_t, a_t)$ 
  - The function  $\delta$  and  $r$  are part of the environment and not necessarily known to the agent

# MDP dependency

- In an MDP the two functions

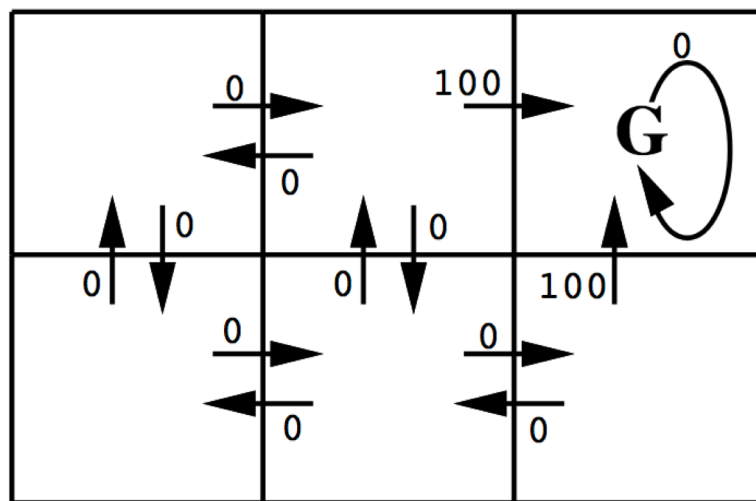
$\delta(s_t, a_t) = s_{t+1}$  which maps the current state to the next state

$r(s_t, a_t) = r_t$  which provides the reward for the action

depend only on the current state and action and not on earlier states and actions

- We consider only case where  $S$  and  $A$  are finite
- We first consider the deterministic case for  $\delta$  and  $r$

# A simple grid world environment



Six grid squares represent six possible states for the agent

Each arrow represents a possible action can take to move to another state

The number with each arrow is immediate reward  $r(s, a)$  agent receives  
 gives reward of 100 for actions entering the goal state G and zero otherwise

G is the goal state for which the only action available is to remain in the state

# How to specify Policy?

- Task of agent is to learn a policy  $\pi: S \rightarrow A$ 
  - For selecting its next action  $a_t$  based on current observed state  $s_t$ , *i.e.*,  $\pi(s_t) = a_t$
- How to specify the policy?
  - One approach: Require policy that produces greatest cumulative reward over time
    - For this we need to define cumulative value  $V^\pi(s_t)$  of a policy  $\pi$  starting from an initial state  $s_t$  as follows

$$\begin{aligned} V^\pi(s_t) &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \\ &= \sum_{i=0}^{\infty} \gamma^i r_{t+i} \end{aligned}$$

Here  $0 \leq \gamma < 1$  is a constant that determines the relative value of delayed vs immediate rewards (see next slide)

- Where sequence of rewards  $r_{t+1}$  is generated by beginning at state  $s_t$  and by repeatedly using policy  $\pi$

# Several Cumulative Reward Definitions

## 1. Discounted Cumulative reward of policy $\pi$ is

$$V^\pi(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{i=0}^{\infty} \gamma^i r_{t+i+1}$$

- If  $\gamma=0$  only the immediate reward is considered
- As we set  $\gamma$  closer to 1, future rewards are given greater emphasis relative to the immediate reward
- Future rewards discounted wrt immediate rewards\*

\*Keynes: "In the long run we are all dead"

## 2. Finite horizon reward

- Undiscounted sum of rewards over finite  $h$  steps

$$V^\pi(s_t) = \sum_{i=0}^h r_{t+i}$$

## 3. Average reward

$$V^\pi(s_t) = \lim_{h \rightarrow \infty} \sum_{i=0}^h r_{t+i}$$



# Agent's Learning Task

- The agent has to learn a policy  $\pi$  that maximizes  $V^\pi(s)$  for all states  $s$

- Where

$$V^\pi(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = \sum_{i=0}^{\infty} \gamma^i r_{t+i+1}$$

- We will call such a policy an optimal policy  $\pi^*$ :

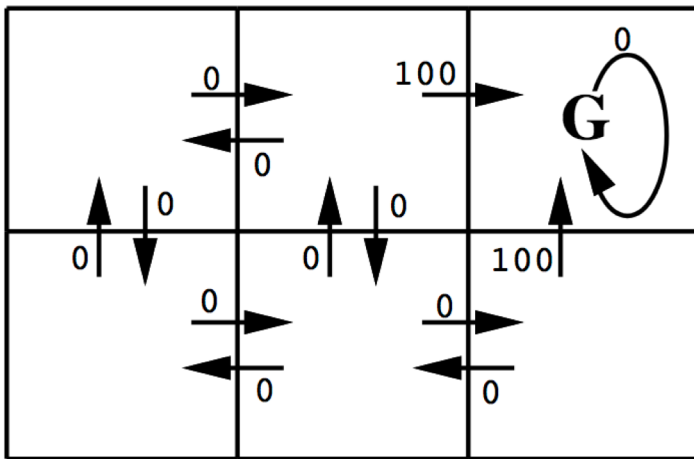
$$\pi^* = \arg \max_{\pi} V^\pi(s), (\forall s)$$

- We denote the value function  $V^{\pi^*}(s)$  by  $V^*(s)$
- It gives the maximum discounted cumulative reward that the agent can obtain starting from state  $s$

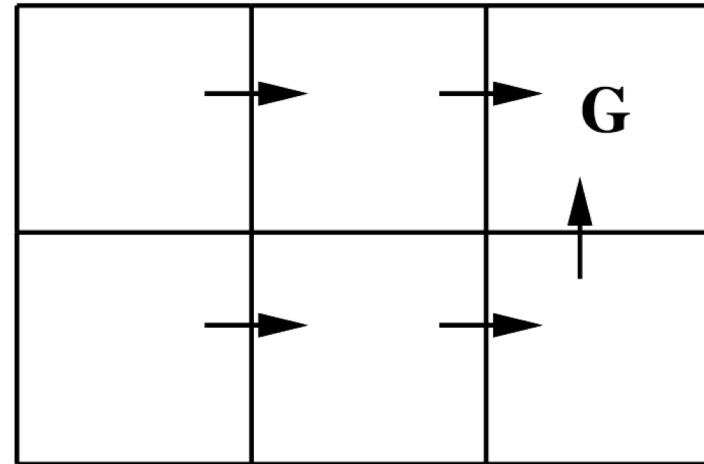
# An optimal policy

- Once the states, actions and immediate rewards are defined, and we have chosen a value for  $\gamma$ , we can determine the optimal policy  $\pi^*$  and its value function  $V^*(s)$

$r(s,a)$  values



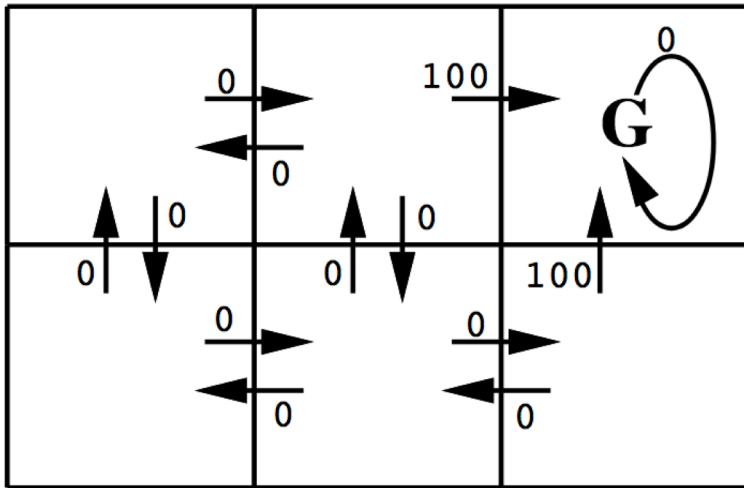
One optimal policy



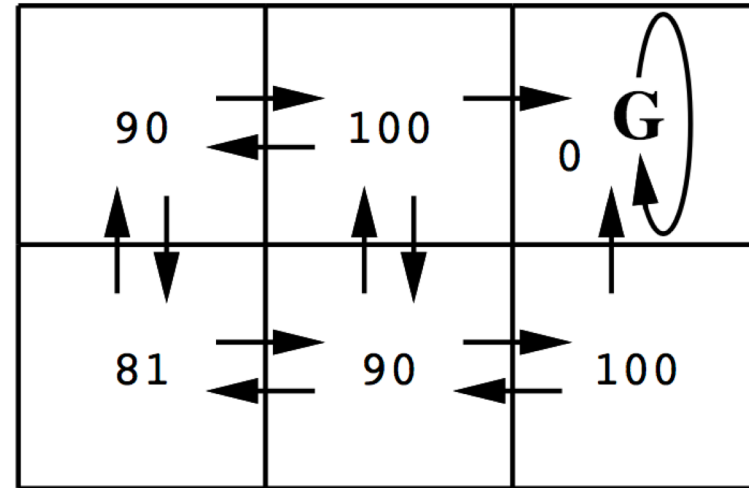
- Optimal policy is obtained by determining values of  $V^*(s)$  and  $Q(s,a)$  from  $r(s,a)$  and  $\gamma=0.9$
  - Selects one action agent will select in any given state
- Optimal policy directs the agent along the shortest path toward the state G

# Value of $V^*$

$r(s, a)$  values



$V^*(s)$  values



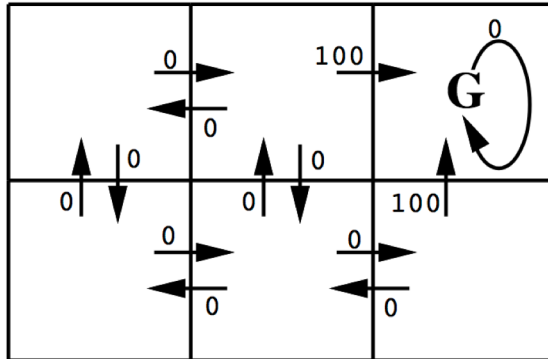
- Value of  $V^*$  for each state shown
- For bottom right state, value is 100 because the optimal policy in this state selects the “move up” action that receives an immediate reward of 100
- For bottom center state, value of  $V^*$  is 90  
This is because the optimal policy will move the agent from this state to the right (generating an immediate reward of zero), then upward (generating an immediate reward of 100)

Thus the discounted future reward is

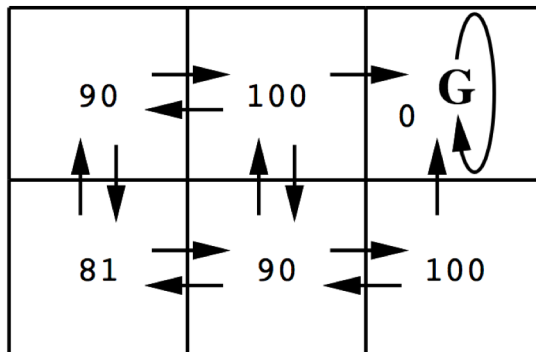
$$0 + \gamma 100 + \gamma^2 0 + \gamma^3 0 + \dots = 90$$

# Q-values

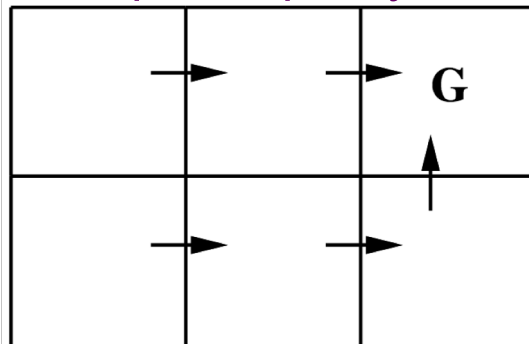
$r(s, a)$  values



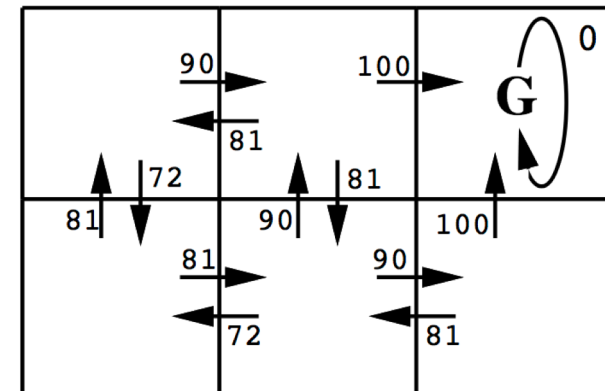
$V^*(s)$  values



Optimal policy



$Q(s, a)$  values



Q values for every state and action are shown. Q-value for each state-action transition is equal to the  $r$  value for this transition plus the  $V^*$  value for the resulting state discounted by  $\gamma$

Optimal policy corresponds to selecting actions with maximal  $Q$  values