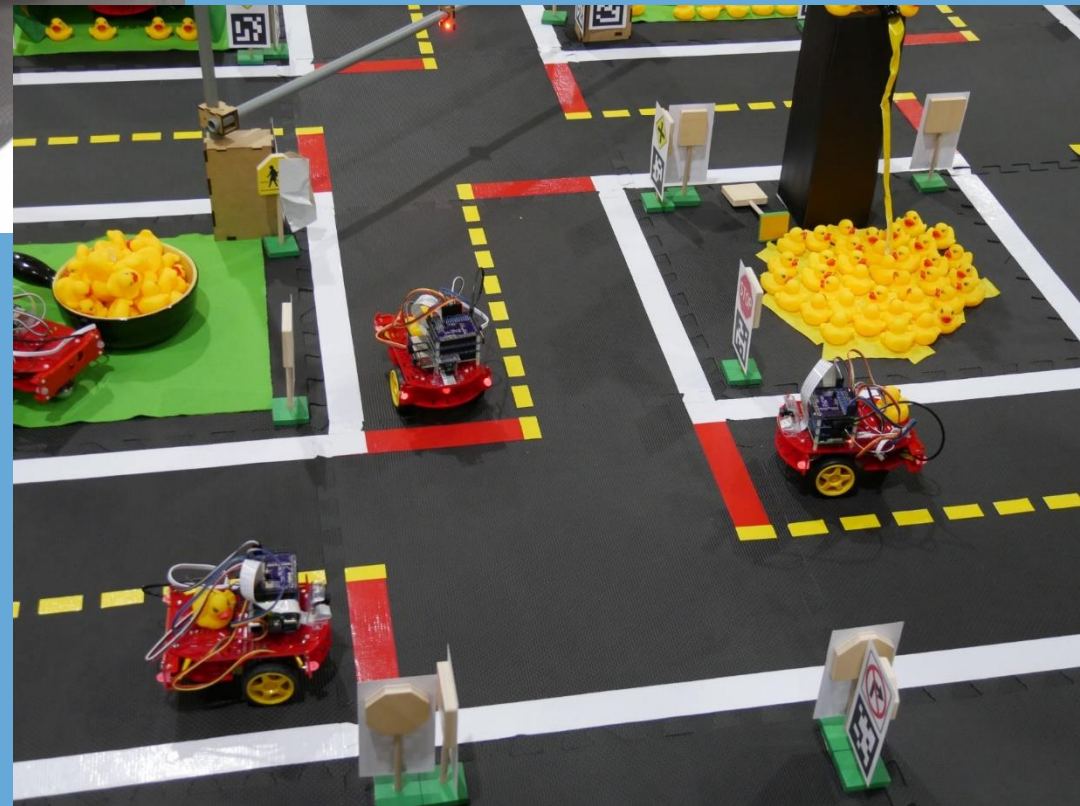


# CS 3630

## Reinforcement Learning: 2



Markov Decision Processes,  
Approaches to  
Reinforcement Learning

# Policies

$$\underline{E[r_h]} = E \left[ \sum_{i=0}^h \gamma^i R(s_i) \mid a_1, \dots, a_h \right]$$

Def: A policy  $\pi: S \rightarrow A$  specifies  $a = \pi(s)$ , the action to take in states.

For a policy  $\pi$

$$\underline{V^\pi(s)} = E[r_\omega(s) \mid \pi]$$

$$= E \left[ \sum_{i=0}^{\infty} \gamma^i R(s_i) \mid s_0 = s, \pi \right]$$

$$= E \left[ \underline{R(s_0)} + \sum_{i=1}^{\infty} \gamma^i R(s_i) \mid s_0 = s, \pi \right]$$

$$= R(s) + E \left[ \gamma \sum_{i=1}^{\infty} \gamma^{i-1} R(s_i) \mid \pi \right]$$

$$\begin{aligned} i &= j+1 \\ j &= i-1 \end{aligned}$$

$$= R(s) + \gamma E \left[ \sum_{j=0}^{\infty} \gamma^j R(s_{j+1}) \mid \pi \right]$$

Expected return under  $\pi$   
from state  $s_{j+1}$

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

Expected return  
for executing  $a = \pi(s)$   
from state  $s$ .

We want the optimal policy,  $\underline{\pi^*}$

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

Def The value function

$$\underline{V^* (= V^{\pi^*})}, V^*: S \rightarrow \mathbb{R}$$

is the value for  $\pi^*$

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s'} T(s, a, s') V^*(s')$$

Expected return  
for executing  
action  $a$  in  
state  $s$

## Bellman Equation

$$\pi^* = \arg \max_a \sum T(s, a, s') V^*(s')$$

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

Bellman  
Eqn

How to find  $V^*$  ??

Suppose we have an estimate  $\underline{V}^k$  of  $V^*$   
and we want to iteratively improve  
s.t.  $V^k \xrightarrow{k \rightarrow \infty} V^*$ , the unique soln  
to Bellman eqn.

Value  
Iteration

$$V^{k+1}(s) = R(s) + \gamma \max_a \sum_{s' \in S'} T(s, a, s') \underline{V}^k(s')$$

Truth

↑  
Best guess at  $V^*$   
at  $k^{\text{th}}$  iteration

If Works

what we would have,  
if  $V^k = V^*$



|       |       |       |              |
|-------|-------|-------|--------------|
| A     | +1    | .3    | -.05         |
| B     | +1    | .3    | 0.1          |
| C     | +1    | .3    | -.25         |
| D     | +1    | .3    | 0.1          |
| E     | +1    | 1.5   | 1.15         |
| F     | +1    | .3    | 0.1          |
| G     | +1    | .3    | -.25         |
| H     | +1    | .3    | 0.1          |
| I     | +1    | -.3   | -.05         |
| J     | +1    | .3    | -.05         |
| K     | +1    | .3    | -.05         |
| L     | +1    | -.3   | -.05         |
| $V^0$ | $V^1$ | $V^2$ | $\gamma=0.5$ |

R ↑  
L ↓

$$V^2(s) = R(s) + 0.5 \max_{a \in \{R, L\}} T(s, a, s') V^1(s')$$

For A, I, J, K, L

$$V^2(s) = -0.2 + 0.5 \max_{a \in \{R, L\}} \{0.25 \times .3 + 0.5 \times .3 + 0.25 \times .3, 0.25 \times 3 + 0.5 \times .3 + 0.25 \times .3\}$$

$$= -0.05$$

S = B

$$V^2(B) = -0.2 + 0.5 \max \{0.3, \underbrace{0.25 \times .3 + 0.5 \times 1.5 + 0.25 \times .3}_{0.6}\}$$

$$= -0.2 + 0.5(0.6)$$

$$= 0.1$$

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

$$V^1(s) = R(s) + 0.5 \sum_a T(s, a, s') V^0(s)$$

For  $s \neq E, R(s) = -.2$

$$V^1(s) = -.2 + 0.5 \max_{R, L} \{0.25 \times 1 + 0.5 \times 1 + 0.25 \times 1, 0.25 \times 1 + 0.5 \times 1 + 0.25 \times 1\}$$

$$= -.2 + .5 \times 1 = 0.3$$

## Policy Iteration

Let  $\pi^i$  be an approximation to  $\pi^*$  at the  $i$ th iteration.

$$\underline{\underline{\pi^{i+1}}}(s) = \arg \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \underbrace{T(s, \pi^i(a), s')}_{\substack{\text{current guess} \\ \text{for } \pi^*}} \underbrace{V^i(s')}_{\substack{\text{current guess} \\ \text{for } V^*}}$$

$$V^{i+1}(s) = R(s) + \gamma \sum_{s' \in \mathcal{S}} \underbrace{T(s, \pi^{i+1}(s), s')}_{\substack{\text{looks like expected return under } \pi^*}} V^i(s')$$

# Reinforcement Learning

We have done a lot of work:

$$\rightarrow \text{MDP} = (S, A, T, R, \gamma)$$

$\rightarrow$  Bellman Eqn  $\rightarrow$  Value Function

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

$\rightarrow$  Optimal Policy

$$\pi^*(s) = \arg \max_a \sum_{s'} T(s, a, s') V^*(s')$$

For Machine Learning

\* We don't know  $T$ .  
 $\hookrightarrow$  world model

\* We don't know  $R(s)$

## R.L. Two Approaches

Passive: The agent has a fixed policy  $\pi$ , and it can learn about  $T$  and/or  $V^\pi$  by executing  $\pi$ .

Active: The policy is not fixed. The agent learns about  $V^*$  and/or  $T$ , while learning to act optimally.

## Direct Utility Estimation

$$V^\pi(s) = E \left[ \sum_{i=0}^{\infty} \gamma^i R(s_i) \mid \pi, s_0 = s \right]$$

Idea execute  $\pi$  from  $s$  many times.  
The average of the returns is a good approximation for  $V^\pi(s)$ .

Can't execute  $\infty$  actions, so execute over a finite horizon, length  $h$ .

$$\begin{aligned} V^\pi(s_0, \dots) &= \sum_{i=0}^{\infty} \gamma^i R(s_i) \\ &= \underbrace{\sum_{i=0}^h \gamma^i R(s_i)}_{\text{Approximation of } V^\pi(s)} + \underbrace{\sum_{i=h+1}^{\infty} \gamma^i R(s_i)}_{\text{Error in approximation}} \end{aligned}$$

$$\begin{aligned} \text{Error} &= \sum_{i=h+1}^{\infty} \gamma^i R(s_i) \leq \sum_{i=h+1}^{\infty} \gamma^i R_{\max} = \sum_{j=0}^{\infty} \gamma^{j+h+1} R_{\max} \\ &= \frac{\gamma^{h+1} R_{\max}}{1-\gamma} \leftarrow \text{Bound on the error} \end{aligned}$$



# Adaptive Dynamic Programming (ADP)

Recall:

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S'} T(s, \pi(s), s') V^\pi(s') \quad (*)$$

Idea at each stage of execution, execute  $a = \pi(s)$  in state  $s$  and arrive to state  $s' \implies \underline{s, a, s'}$ , we observe  $R(s')$

At each stage

- For each  $\bar{s}$ , update  $\hat{T}(s, a, \bar{s}) = \frac{N_{sa\bar{s}}}{N_{sa}}$    
 *new estimate of  $T$*
- Update  $\hat{V}^\pi(s) \leftarrow R(s) + \gamma \sum_{s'} \hat{T}(s, a, s') \hat{V}^\pi(s')$    
 *new estimate for  $V^\pi(s)$*    
 *prev. estimate of  $V^\pi$*

$N_{sa\bar{s}}$  = # of times we experience  $s, a, \bar{s}$

$N_{sa}$  = # of times we execute action  $a$  from state  $s$

$\hookrightarrow$  Keep in a table

|| USE D.P. to implement this



Temporal Difference Learning (TD Learning) ← Model-Free  
Execute a in states s and reach s'. ⇒ s, a, s' Approach.

$$R(s) + \gamma \pi' \quad R(s) + \gamma V^\pi(s') \rightsquigarrow \text{after experience } s, a, s'$$

Before experiencing s, a, s', the best guess for return is  $V^\pi(s)$

$$\underbrace{[R(s) + \gamma V^\pi(s')] - V^\pi(s)}_{\text{Temporal Difference.}}$$

No Model of T

Updating Scheme:

$$\underbrace{V^\pi(s)}_{\text{TD equation}} \leftarrow V^\pi(s) + \underbrace{\alpha}_{\text{Learning rate}} [R(s) + \gamma V^\pi(s) - V^\pi(s)]$$

ADP VS. TD

ADP:  $V^\pi$  is made to agree with all past experience.

TD:  $V^\pi$  is made to agree only with current experience.

## Active Learning

- If we have a model,  $T$ , we can use  $T$  to build an approximation of  $V^*$ , and thus  $\pi^*$ . [e.g., using Value or Policy Iteration]

Q Given  $\hat{T} \Rightarrow \hat{V}^*, \hat{\pi}^*$  should we

a. execute  $\hat{\pi}^*$  (s)  $\leftarrow$  Exploitation

b. Try some other action,  $a$ ,  $\leftarrow$  Exploration  
maybe finding better  $\hat{V}^*$ .

$\rightarrow$  Key problem: Balancing Exploration vs. Exploitation

Many available Algorithms to do this---

Q-Learning: For active learning, we need to explicitly consider the action  $a$  when deciding what to do.

$V^*(s)$  does not explicitly consider the action to be performed.

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

*action appears.*

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

$Q(a, s) = Q$  function

$$\bullet \underline{V^*(s)} = \max_a Q(a, s)$$

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

constraint eqn  
for  $Q$  function

$$Q(a, s) = R(s) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(a', s')$$



# TD Q-Learning

Given  $Q$  at some stage of Learning.  
We execute action  $a$  in states, arrive to  $s'$   
 $\Rightarrow \underline{s, a, s'}$

$$Q(a, s) \leftarrow \underbrace{Q(a, s)}_{\text{current estimate of } Q} + \alpha \left[ \underbrace{R(s) + \gamma \max_{a'} Q(a', s')}_{\text{estimate of } Q \text{ based on having seen } s, a, s'} - \underbrace{Q(a, s)}_{\text{estimate of } Q \text{ before seeing } s, a, s'} \right]$$

Learning rate.

Note We never build a model for  $T \Rightarrow$  Model-Free.

END