

# CS 3630

## Reinforcement Learning



# Reinforcement Learning

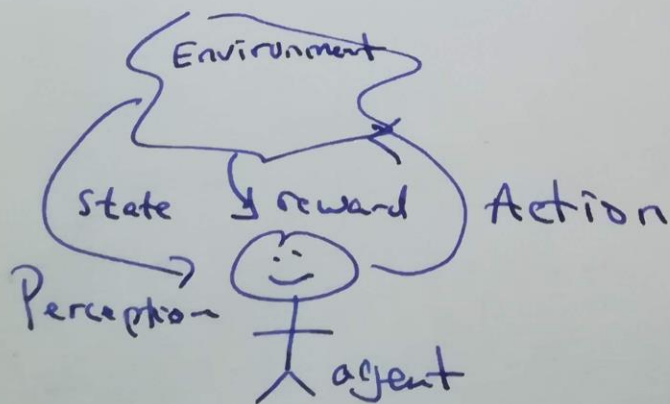
## "Typical" Machine Learning

- Lots of Data
- passive Learner
- Learning = pattern analysis  
function approx.

## Examples

- \* Find photos of cats.
- \* product recommendation
- \* Text Translation

## Robots are Not Like This!



## RL:

RL is a process of modifying behavior by rewarding desired outcomes.

- Dogs - Treats, learn tricks  
- punishment  
→ negative reward
- Kids - 100/A on report
- students: grades  
joy, satisfaction  
Starting Salary

1. Robot senses world
2. Robot decides & executes action
3. World state changes
4. Robot Receives a reward. }  
Robot senses world
5. Robot update its "strategy"
6. Go to 2



## Questions:

- Mathematical Model
  - states
  - actions
  - uncertainty (!)
- Mathematical formalism for reward
- Given the above  
How to compute an optimal strategy/policy
- Now... suppose we don't know the models for world, actions... how can we learn them?

- 
- Markov Process
  - Markov Decision Process
  - Value function & how to compute it
  - R.L.

## Markov Processes

### Simplest case:

- discrete time,  $k=0, 1, 2, \dots$
- discrete states,  $S$  set of states.
- A set of transition probabilities

$$T: S \times S \rightarrow [0, 1]$$

probability

$$P \{S_{k+1} | S_k\} = \underline{T}(s, s')$$
$$= T(S_k, S_{k+1})$$

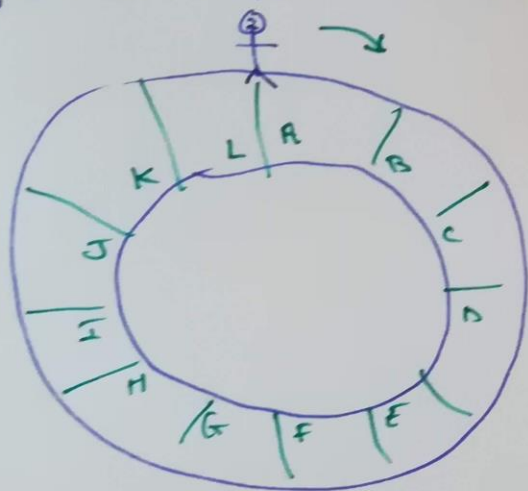
$$T(s, s') = \text{prob} \{ \text{arriving to state } s', \\ \text{given we are now in} \\ \text{state } s \}$$

---

↳ M.P. evolves "autonomously"

### Example:

Sisyphus has the job of, every day, throwing a large rock, on a circular track.



$L_k$  = distance of throw on  $k^{\text{th}}$  day

$$P \{ L_k = 1 \} = 0.25$$

$$P \{ L_k = 2 \} = 0.5$$

$$P \{ L_k = 3 \} = 0.25$$

$$S = \{A, B, \dots, L\}$$

$$T(A, B) = 0.25 \quad T(A, C) = 0.5 \quad T(A, D) = 0.25$$

$\vdots$

$$T(J, K) = 0.25 \quad T(J, L) = 0.5 \quad T(J, A) = 0.25$$

$$T(K, L) = 0.25 \quad T(K, A) = 0.5 \quad T(K, B) = 0.25$$

$$T(S, S') = 0 \text{ otherwise}$$

### Properties:

- Stationary:  $T$  does not change over time

$$P \{ S_{k+1} | S_0, S_1, \dots, S_k \} = P \{ S_{k+1} | S_k \}$$

### Markov Property

$$P \{ S_3 = E | S_0 = A, S_1 = C, S_2 = D \} = P \{ S_3 = E | S_2 = D \}$$

$$P \{ S_3 = E | S_0 = A, S_1 = B, S_2 = D \} = \uparrow$$





## Expectation

Suppose a random variable  $X$  takes values from the set  $\{c_1, c_2, \dots, c_n\}$ . The expected value of  $X$  is defined as

$$E[X] = \sum_{i=1}^n P\{X=c_i\} \times c_i$$

Example: roll one die,  $X$  shows

$$E[X] = \sum_{i=1}^6 \frac{1}{6} \times i = \underline{3.5}$$

↓  
all values equally likely.

Intuition: Perform this many times

Average of  $X$ 's  $\rightarrow E[X]$

## Generalization

$$E[X_1 + X_2] = \sum_i \sum_j P\{X_1=c_i, X_2=c_j\} (c_i + c_j)$$

Two dice

$$E[X_1 + X_2] = \sum_i \sum_j \frac{1}{36} (i+j) = 7$$

$$E[\sum_i X_i] = \sum_i P\{X_i=c_i\} \sum_j (c_j)$$

For Sisyphus:  $E[R(S_0) + R(S_1) + R(S_2)]$   
given  $a_1 = L, a_2 = L$ ??

We know

$$\left\{ \begin{array}{l} R(S_0) = -0.2 \text{ because } S_0 = A \\ R(S_1) = -0.2, S_1 \in \{B, C, D\} \\ R(S_2) = \begin{cases} +1 & S_2 = E \\ -0.2 & \text{Else} \end{cases} \end{array} \right.$$

## Expected return

Define return  $r_h = \sum_{i=0}^h R(s_i)$ . ←  $\otimes$

$r_2$  for sisyphus =  $R(s_0) + R(s_1) + R(s_2)$

$$E[r_2] = E[R(s_0) + R(s_1) + R(s_2)]$$

$$= P\{S_2 = E\} \times 0.6 + P\{S_2 \neq E\} \times -0.6 \otimes$$

because  $r_h$  has only two possible values

$$-0.2 + -0.2 + -0.2 = -0.6 \leftarrow S_2 \neq E$$

$$-0.2 + -0.2 + 1 = 0.6 \leftarrow S_2 = E$$

To arrive  $S_2 = E$

$S_0$	$S_1$	$S_2$	Probability
A	B	E	$(0.25) \times (0.25) = 0.0625$
A	C	E	$(0.5) \times (0.5) = 0.25$
A	D	E	$(0.25) \times (0.25) = 0.0625$

~~0.375~~

$$P(S_2 = E) = 0.375$$

$$P\{S_3 \neq E\} = 1 - P\{S_3 = E\} \\ = 0.625$$

This is great for finite horizons (i.e. only consider a finite # of stages).

Suppose  $h \rightarrow \infty$  ??

To deal with this, use discounted

rewards · discount factor

$$r_h = \sum_{i=0}^h \gamma^i R(s_i)$$

for  $0 < \gamma < 1$ .

$$\lim_{h \rightarrow \infty} \sum_{i=0}^h \gamma^i R(s_i) \leq \sum_{i=0}^{\infty} \gamma^i R_{\max}$$

$$\sum_{i=0}^{\infty} \gamma^i = \frac{1}{1-\gamma} \quad 0 < \gamma < 1$$

$$\sum_{i=0}^{\infty} \gamma^i R(s_i) \leq \frac{R_{\max}}{1-\gamma}$$

## 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  $\pi$

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

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

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

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

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

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

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

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.

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

$$V^{i+1}(s) = R(s) + \gamma \sum_{s' \in \mathcal{S}} T(s, \underline{\underline{\pi^{i+1}}}(s), s') V^i(s')$$

looks like expected return under  $\pi$

# 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  $h$ .

$$V^\infty(s_0, \dots) = \sum_{i=0}^{\infty} \gamma^i R(s_i)$$

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

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

# 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  $s$

↳ 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)  $\longleftarrow$  Exploitation

b. Try some other action,  $a$ ,  $\longleftarrow$  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

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

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

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