

## Section 2

Function approximation

Supervised learning  $y = f(x)$ , find the

function to approximate to  $x \rightarrow y$

Unsupervised learning:

Clustering

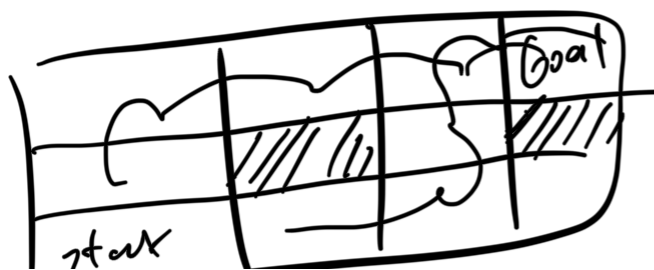
Using  $x$ , find function

Reinforcement learning:

We're going to be given  $X$  and  $Z$ ,

find  $f$  that generates  $y$ 's.

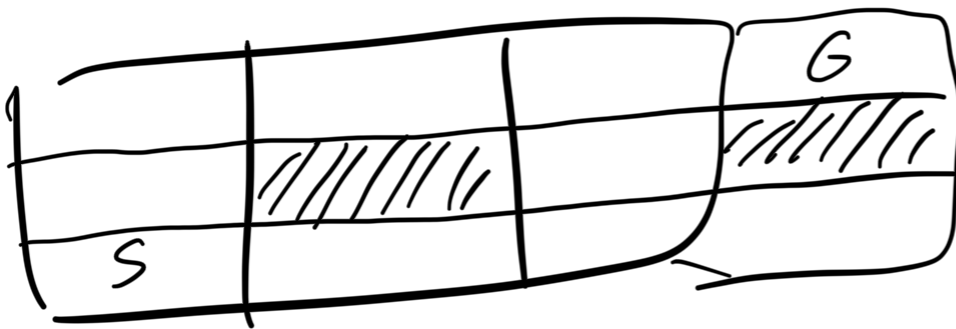
---



Shortest Path

RRUUR

V V R R R



Actions executes (0.8)

Move at angle to left or right (0.2)

Reliability for V V R R R

.32776

$$.8^5 = .32768$$

$$.1^4 \cdot .8 = .00008$$

$$4 \text{ of them } \frac{+}{=} .32776$$

but 1 goes right

Markov Decision Process

Markov property:

States:  $S$

→ transition function

only present matters  
rules are stationary

Model  $T(s, a, s') \sim P_r(s' | s, a)$

Action:  $A(s), A$

Reward:  $R(s), R(s, a), R(s, a, s')$

Policy:  $\pi(s) \rightarrow a$ ,  $\pi^*$   
↳ (canonical, order)

↓ optimal  
max reward

$\langle s, a, r \rangle$

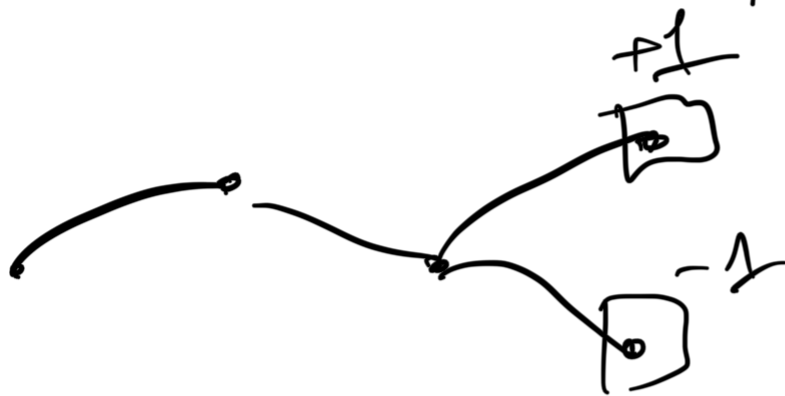
$R: \mathbb{Z}$  for  $y: f(x, r)$

$f: \pi$

$y: a$

$x: s$

delegated server, inner layers matter



Typical  
Credit  
Assignment  
Problem

→	→	→	0
↑	///	↑	///
↑	←	←	←

to avoid possibly pinning to 1.

$R(s): -.04$

(except for peak +1 or tail -1)

Quiz

$R(s): +2$

←  
←  
↑

	?	+1
///	?	-1
	?	?

$$R(s) = -2$$

$\begin{matrix} \leftarrow \rightarrow \\ \downarrow \\ \rightarrow \\ \rightarrow \\ \rightarrow \uparrow \end{matrix}$

		2	+1
		?	-1
		?	?
		?	?

## Sequence of Rewards

- Infinite Horizons
- Utility of sequences

$$U(s_0, s_1, \dots, s_n) = \sum_{t=0}^{\infty} R(s_t)$$

Quiz:

+1 +1 +1 +1 +1 +1  
 ~~~~~ ...  $\infty$

~~~~~ ...  
 +1 +1 +2 +1 +2 ...  $\infty$

neither is better. Sum<sup>d</sup> utility is  $\infty$   
for both.

if our rewards are always +1,  
then our moment doesn't really matter.

$$\sum_{t=0}^{\infty} \gamma^t \cdot R(S_t) \quad 0 \leq \gamma < 1$$

rewards put smaller over time. if it's  
zero, only

$$\frac{R_{\max}}{1-\gamma} = \sum_{t=0}^{\infty} \gamma^t \cdot R_{\max} \quad \text{immediate rewards}$$

↓  
discounted sums.

discounted  $\rightarrow$  generic  
infinite  $\rightarrow$  finite

$$(\sum \gamma^t) R_{\max}$$

$$x: (\gamma^0 + \gamma^1 + \gamma^2 + \gamma^3 \dots)$$

$$x: \gamma^0 + \gamma \cdot x$$

$$x: \frac{1}{1-\gamma} \cdot R_{\max}$$

Policies:  $\arg\max_x: E \left[ \sum_t \gamma^t R(s_t) \mid \pi \right]$

$\pi^*$ :  $\downarrow$   
expected value

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

$Q(s)^*$

Rein: immediate feedback

n. short

Utility: long-term future

(Reverse now, and reveal of future)

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

$U(s) \equiv U^*(s) \rightarrow$  the utility of a state.

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

## Bellman Equation

N of states

$n$  equations,  $n$  unknowns.

max is non-linear

- Start with arbitrary utilities

[illegible]



- update based on neighbors.  
 ↪ repeat until convergence

We use all values from  $U$ 's from previous to update current.

$$\hat{U}(s)_{t+1} = \underbrace{R(s)}_{\rightarrow \text{reward}} + \gamma \cdot \max_a \cdot \sum_{s'} T(s, a, s') \cdot \hat{U}_t(s')$$

Value Iteration.

Quiz:

.36

|   |     |       |    |
|---|-----|-------|----|
|   |     | X     | G  |
|   | /// | -0.04 | -1 |
| S |     |       |    |

$$\gamma = \frac{1}{2}, R(s) = -0.04, U_0(s) = 0$$

$$U_1(x) = 0$$

$$\cancel{U_1^*} - .04 + \overset{\text{we go to the right}}{\frac{1}{2} \cdot [0 + 0 + 0.8]} = \underline{\underline{.36}}$$

$$\cancel{U_2^*} - 0.04 + \frac{1}{2} \left( \begin{array}{c} 0.36 + 0.8 \\ - 0.04 \end{array} \right) = \underline{\underline{0.326}}$$

finding policies:

Start with  $\pi_0$

evaluate given  $\pi_t$  calculate  $V_t = U_t^{\pi_t}$

...

... ..

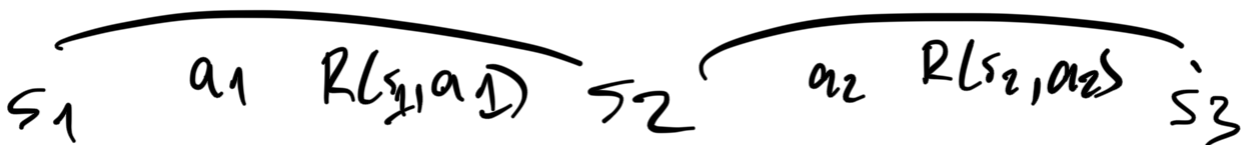
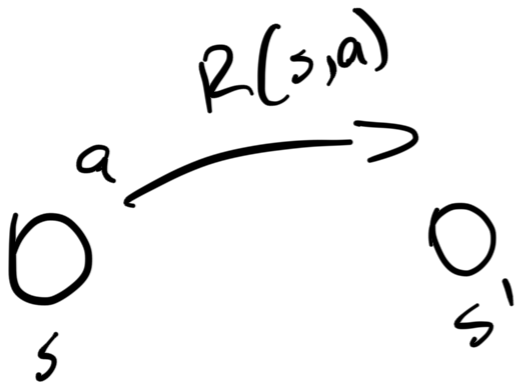
$$\text{Improve } \pi_{t+1} : \arg \max_a \sum_{s'} T(s, a, s') \cdot U_t(s)$$

$$U_t(s) = R(s) + \gamma \sum_{s'} T(s, \pi_t(s), s') \cdot U_t(s')$$

↗ This is linear unlike max eq.

Policy iteration

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



$$V(s_1) : \max_{a_1} \left[ R(s_1, a_1) + \gamma \cdot \sum_{s_2} T(s_1, a_1, s_2) \cdot \right.$$

$$\downarrow \quad \max_{a_2} \left[ R(s_2, a_2) + \gamma \cdot \sum_{s_3} T(s_2, a_2, s_3) \dots \right]$$

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

$$\max_{a'} Q(s', a')$$

|   | V  | Q  | C  |
|---|--|--|--|
| V | $V(s) : V(s)$  | $V(s) : \max_a Q(s, a)$                                      | $V(s) : \max_a (R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s'))$ |
| Q | $Q(s, a) : R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s')$ | $Q(s, a) : Q(s, a)$  | $Q(s, a) : R(s, a) + \gamma \sum_{s'} T(s, a, s') V(s')$       |
| C | $C(s, a) : \gamma \sum_{s'} T(s, a, s') V(s') + V(s)$    | $C(s, a) : \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$ | $C(s, a) : C(s, a)$  |

72