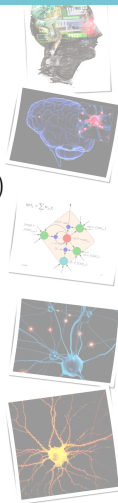


An introduction to Reinforcement Learning (with Neural Networks and Causality)

Spyros Samothrakis
Research Fellow, IADS
University of Essex

November 14, 2016



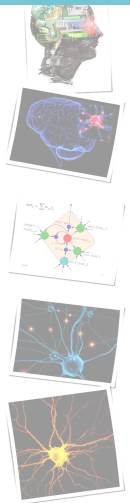
1 / 70

Introduction & Motivation

Markov Decision Process (MDPs)

Planning

Model Free Reinforcement Learning



2 / 70

WHAT IS REINFORCEMENT LEARNING?

- *Reinforcement learning is the study of how animals and artificial systems can learn to optimize their behavior in the face of rewards and punishments* – Peter Dyan, Encyclopedia of Cognitive Science
- **Not** supervised learning - the animal/agent is not provided with examples of optimal behaviour, it has to be discovered!
- **Not** unsupervised learning either - we have more guidance than just observations

3 / 70

LINKS TO OTHER FIELDS

- It subsumes most artificial intelligence problems
- Forms the basis of most modern intelligent agent frameworks
- Ideas drawn from a wide range of contexts, including psychology (e.g., Skinner's "Operant Conditioning"), philosophy, neuroscience, operations research, **Cybernetics**
- Modern Reinforcement Learning research has fused with Neural Networks research

4 / 70

EXAMPLES OF REINFORCEMENT LEARNING CLOSER TO CS

- Play backgammon/chess/go/poker/any game (at human or superhuman level)
- Helicopter control
- Learn how to walk/crawl/swim/cycle
- Elevator scheduling
- Optimising a petroleum refinery
- Optimal drug dosage
- Create NPCs

5 / 70

THE MARKOV DECISION PROCESS

- The primary abstraction we are going to work with is the Markov Decision Process (MDP).
- MDPs capture the dynamics of a mini-world/universe/environment
- An MDP is defined as a tuple $\langle S, A, T, R, \gamma \rangle$ where:
 - S , $s \in S$ is a set of states
 - A , $a \in A$ is a set of actions
 - $R : S \times A$, $R(s, a)$ is a function that maps state-actions to rewards
 - $T : S \times S \times A$, with $T(s'|s, a)$ being the probability of an agent landing from state s to state s' after taking a
 - γ is a discount factor - the impact of time on rewards

6 / 70

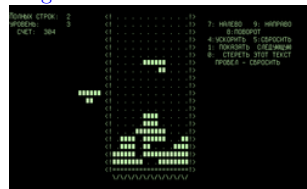
THE MARKOV PROPERTY AND STATES

- ▶ States represent sufficient statistics.
- ▶ Markov Property ensures that we only care about the present in order to act - we can safely ignore past states
- ▶ Think Tetris - all information can be captured by a single screen-shot

First DOS Version



Original Tetris



7 / 70

AGENTS, ACTIONS AND TRANSITIONS

- ▶ An agent is an entity capable of actions
- ▶ An MDP can capture any environment that is inhabited either by
 - ▶ Exactly one agent
 - ▶ Multiple agents, but only one is adaptive
- ▶ Notice how actions are part of the MDP - notice also how the MDP is a “world model”
- ▶ The agent is just a “brain in a vat”
- ▶ The agent perceives states/rewards and outputs actions
- ▶ Transitions specify the effects of actions in the world (e.g., in Tetris, you push a button, the block spins)

8 / 70

MORE ON STATES, AGENTS AND ACTIONS

- ▶ Pick a game
- ▶ What would be state in the game?
 - ▶ Do agents/NPCs have access to it?
- ▶ Do agents/NPCs have access to actions
- ▶ Do agents/NPCs have access to transitions?
- ▶ We will come back to these questions later

9 / 70

REWARDS AND THE DISCOUNT FACTOR

- ▶ Rewards describe state preferences
- ▶ Agent is happier in some states of the MDP (e.g., in Tetris when the block level is low, a fish in water, pacman with a high score)
- ▶ Punishment is just low/negative reward (e.g., being eaten in pacman)
- ▶ γ , the discount factor,
 - ▶ Describes the impact of time on rewards
 - ▶ “I want it now”, the lower γ is the less important future rewards are
- ▶ There are no “springs/wells of rewards” in the real world
 - ▶ What is “human nature”?

10 / 70

EXAMPLES OF REWARD SCHEMES

- ▶ Scoring in most video games
- ▶ The distance a robot walked for a bipedal robot
- ▶ The amount of food an animal eats
- ▶ Money in modern societies
- ▶ Army medals (“Gamification”)
- ▶ Vehicle routing
 - ▶ (-)Fuel spent on a flight
 - ▶ (+) Distance Covered
- ▶ Cold/Hot
- ▶ Do you think there is an almost universal reward in modern societies?

11 / 70

LONG TERM THINKING

- ▶ It might be better to delay satisfaction
- ▶ Immediate reward is not always the maximum reward
- ▶ In some settings there are no immediate rewards at all (e.g., most solitaire games)
- ▶ MDPs and RL capture this
- ▶ “Not going out tonight, study”
- ▶ Long term investment

12 / 70

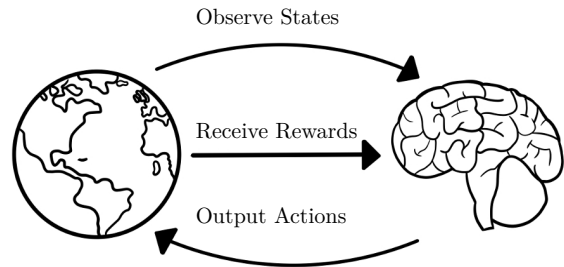
POLICY

- ▶ The MDP (the world) is populated by an agent (an actor)
- ▶ You can take actions (e.g., move around, move blocks)
- ▶ The type of actions you take under a state is called the *policy*
- ▶ $\pi : S \times A, \pi(s, a) = P(a|s)$, a probabilistic mapping between states and actions
- ▶ Finding an optimal policy is *mostly* what the RL problem is all about

13 / 70

THE FULL LOOP

- ▶ See how the universe described by the MDP defines actions, not just states and transitions
- ▶ An agent needs to act upon what it perceives
- ▶ Notice the lack of body - “brain in a vat”. Body is assumed to be part of the world.



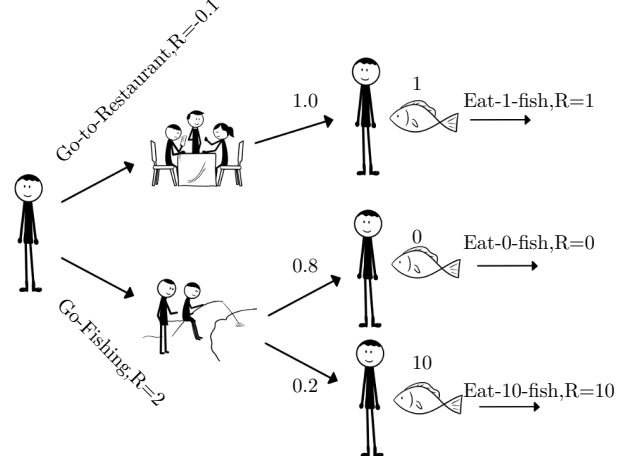
14 / 70

FISHING TOON

- ▶ Assume a non-player character (let's call her *toon*)
- ▶ Toon is Hungry!
- ▶ Eating food is rewarding
- ▶ Has to choose between going fishing or going to the restaurant (to eat fish)
 - ▶ Fishing can get you better quality of fish (more reward), but you might also get no fish at all (no reward)!
 - ▶ Going to the restaurant is a low-risk, low-reward alternative

15 / 70

FISHING TOON: PICTORIAL DEPICTION



16 / 70

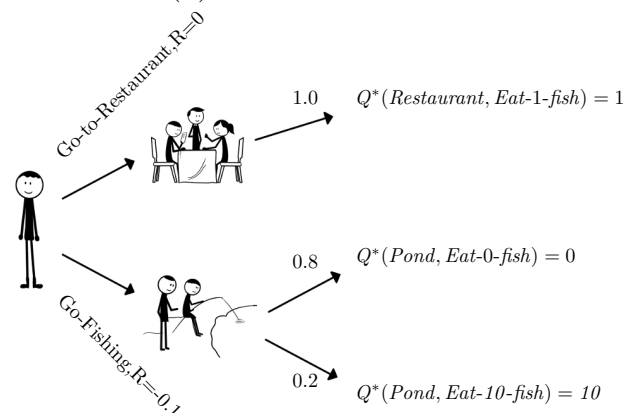
SUM OF EXPECTED REWARDS

- ▶ Our toon has to choose between two different actions
- ▶ Go-To-Restaurant or Go-Fishing
- ▶ We assume that toon is interested in maximising the *expected sum* of happiness/reward
- ▶ Let's first see what happens if we start with a random policy

Policy	Policy Value	Q-Values
$\pi(\text{Start}, \text{Go-Fishing})$	0.5	
$\pi(\text{Start}, \text{Go-to-Restaurant})$	0.5	
$\pi(\text{Restaurant}, \text{Eat-1-fish})$	1	
$\pi(\text{Pond}, \text{Eat-0-fish})$	1	
$\pi(\text{Pond}, \text{Eat-10-fish})$	1	

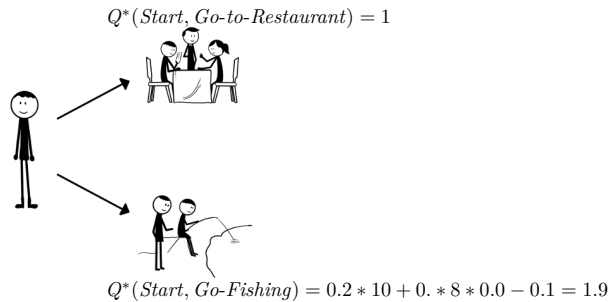
17 / 70

RANDOM POLICY (1)



18 / 70

RANDOM POLICY (2)



19 / 70

TABLE

Policy	Policy Value	Q-Values
$\pi(Start, Go-Fishing)$	0.5	1
$\pi(Start, Go-to-Restaurant)$	0.5	1.9
$\pi(Restaurant, Eat-1-fish)$	1	1
$\pi(Pond, Eat-0-fish)$	1	0
$\pi(Pond, Eat-10-fish)$	1	10

The V-Value of state *Start* is $V(Start) = 0.5 * 1 + 0.5 * 1.9 = 1.45$

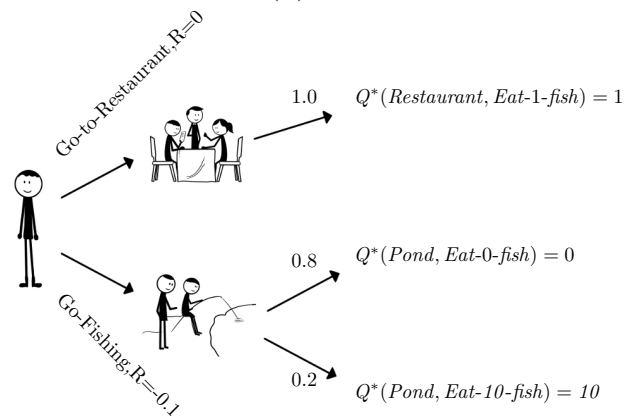
20 / 70

WHAT IF WE ARE ASKED TO FIND OUT THE OPTIMAL POLICY?

Policy	Policy Value	Q-Values
$\pi(Start, Go-Fishing)$?	1
$\pi(Start, Go-to-Restaurant)$?	1.9
$\pi(Restaurant, Eat-1-fish)$	1	1
$\pi(Pond, Eat-0-fish)$	1	0
$\pi(Pond, Eat-10-fish)$	1	10

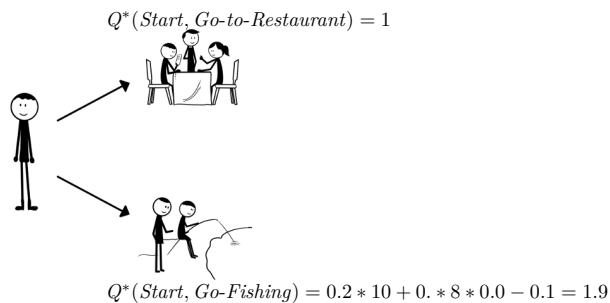
21 / 70

REASONING BACKWARDS (1)



22 / 70

REASONING BACKWARDS (2)



23 / 70

TABLE

Policy	Policy Value	Q-Values
$\pi(Start, Go-Fishing)$	0	1
$\pi(Start, Go-to-Restaurant)$	1	1.9
$\pi(Restaurant, Eat-1-fish)$	1	1
$\pi(Pond, Eat-0-fish)$	1	0
$\pi(Pond, Eat-10-fish)$	1	10

The V-Value of state *Start* is $V^*(Start) = \max\{1, 1.9\} = 1.9$

24 / 70

INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 275 334 302">CORRECT ACTION</h2> <ul data-bbox="144 367 602 472" style="list-style-type: none"> ▶ Toon should go Go-Fishing ▶ Would you do the same? ▶ Would a pessimist toon do the same? ▶ We just went through the following equation: $Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s' s, a) \max_{a' \in A} Q^*(s', a')$ <ul data-bbox="144 564 732 672" style="list-style-type: none"> ▶ Looks intimidating - but it's really simple ▶ Let's have a look at another example <ul data-bbox="190 625 732 672" style="list-style-type: none"> ▶ How about toon goes to the restaurant after failing to fish? ▶ How would that change the reward structure? <p data-bbox="753 753 797 768">25 / 70</p>	<h2 data-bbox="829 275 1016 302">AGENT GOALS</h2> <ul data-bbox="875 359 1498 684" style="list-style-type: none"> ▶ The agent's goal is to maximise its long term reward $\mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s^t, a^t) \right]$ ▶ Risk Neutral Agent - think of the example above ▶ Rewards can be anything, but most agents receive rewards only in a very limited amount of states (e.g., fish in water) ▶ What if your reward signal is only money? <ul data-bbox="920 564 1498 684" style="list-style-type: none"> ▶ Sociopathic, egotistic, greed-is-good Gordon Gekko (<i>Wall Street</i>, 1987) ▶ No concept of “externalities” - agents might wreak havoc for marginal reward gains ▶ Same applies to all “compulsive agents” - think Chess <p data-bbox="1484 753 1528 768">26 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 827 529 854">SEARCHING FOR A GOOD POLICY</h2> <ul data-bbox="144 974 764 1136" style="list-style-type: none"> ▶ One can possibly search through all combinations of policies until she finds the best ▶ Slow, does not work in larger MDPs ▶ Exploration/Exploitation dilemma <ul data-bbox="190 1089 764 1136" style="list-style-type: none"> ▶ How much time/effort should be spend exploring for solutions? ▶ How much time should be spend exploiting good solutions? <p data-bbox="753 1304 797 1318">27 / 70</p>	<h2 data-bbox="829 827 959 854">PLANNING</h2> <ul data-bbox="875 932 1498 1199" style="list-style-type: none"> ▶ An agent has access to model, i.e. has a copy of the MDP (the outside world) in its mind ▶ Using that copy, it tries to “think” what is the best route of action ▶ It then executes this policy on the real world MDP ▶ You can't really copy the world inside your head, but you can copy the dynamics ▶ “This and that will happen if I push the chair” ▶ Thinking, introspection... ▶ If the model is learned, sometimes it's called “Model Based RL” <p data-bbox="1484 1304 1528 1318">28 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 1373 732 1436">BELLMAN EXPECTATION EQUATIONS / BELLMAN BACKUPS</h2> <ul data-bbox="144 1520 732 1745" style="list-style-type: none"> ▶ The two most important equations related to MDP ▶ Recursive definitions ▶ $V^{\pi}(s) = \sum_{a \in A} \pi(s, a) \left(R(s, a) + \gamma \sum_{s' \in S} T(s' s, a) V^{\pi}(s') \right)$ ▶ $Q^{\pi}(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s' s, a) \sum_{a' \in A} \pi(s', a') Q^{\pi}(s', a')$ ▶ Called V-Value(s) (state-value function) and Q-Value(s) (state-action value function) respectively ▶ Both calculate the expected rewards under a certain policy <p data-bbox="753 1850 797 1864">29 / 70</p>	<h2 data-bbox="829 1373 1187 1400">LINK BETWEEN V^{π} AND Q^{π}</h2> <ul data-bbox="875 1535 1398 1667" style="list-style-type: none"> ▶ V and Q are interrelated ▶ $V^{\pi}(s) = \sum_{a \in A} \pi(s, a) Q^{\pi}(s, a)$ ▶ $Q^{\pi}(s, a) = R(s, a) + \sum_{s' \in S} T(s' s, a) V^{\pi}(s')$ ▶ V-values are defined on states, Q-values on policies! <p data-bbox="1484 1850 1528 1864">30 / 70</p>

OPTIMAL POLICY AND THE BELLMAN OPTIMALITY EQUATION

- ▶ An optimal policy can be defined in terms of Q-values
- ▶ It is the policy that maximises Q values
- ▶ $V^*(s) = \max_{a \in A} R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) V^*(s')$
- ▶ $Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) \max_{a' \in A} Q^*(s', a')$
- ▶ $\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} Q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$

31 / 70

LINK BETWEEN V^* AND Q^*

- ▶ Again, they are interrelated
- ▶ $V(s)^* = \max_{a \in A} Q^*(s, a)$
- ▶ $Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) V^*(s')$
- ▶ Let's assume that toon has another option
- ▶ She can go and buy and eat some meat with a reward of 1.5
- ▶ Or go down the fish route
- ▶ Write down the MDP
 - ▶ Find out the new Q and V values with:
 - ▶ Toon acting randomly on choosing a decision point
 - ▶ Toon choosing action *Go-Fishing*
 - ▶ Toon choosing action *Go-to-Restaurant*

32 / 70

AGENTS REVISITED

- ▶ An Agent can be composed of a number of things
- ▶ A policy
- ▶ A Q-Value/and or V-Value Function
- ▶ A Model of the environment (the MDP)
- ▶ Inference/Learning Mechanisms
- ▶ ...
- ▶ An agent has to be able to *discover a policy* either on the fly or using Q-Values
- ▶ The Model/Q/V-Values serve as intermediate points towards constructing a policy
- ▶ Not all RL algorithms use that (but most do)...

33 / 70

SIMPLIFYING ASSUMPTIONS

- ▶ Assume deterministic transitions
- ▶ Thus, taking an action on a state will lead only to ONE other possible state for some action a_c
 - ▶ $T(s'|s, a_i) = \begin{cases} 1 & \text{if } a_i = a_c \\ 0 & \text{otherwise} \end{cases}$
 - ▶ $V^*(s) = \max_{a \in A} [R(s, a) + \gamma V^*(s')]$
 - ▶ $Q^*(s, a) = R(s, a) + \gamma \max_{a' \in A} Q(s', a')$
- ▶ It is easier now to solve for problems that have loops in them
- ▶ We can also attempt to learn Q-Values without a model!
- ▶ All we need in order to find the optimal policy is $Q(s, a)$

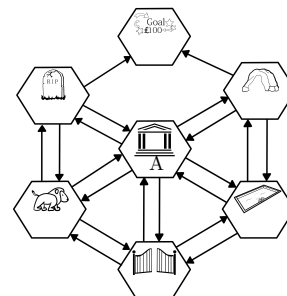
34 / 70

DETERMINISTIC Q-LEARNING (1)

- ▶ The policy is deterministic from start to finish
- ▶ We will use $\pi(s) = \arg \max_{a \in A} Q(s, a)$ to denote the optimal policy
- ▶ The algorithm now is:
 - ▶ Initialise all $Q(s, a)$ to low values
 - ▶ Repeat:
 - ▶ Select an action a using an exploration policy
 - ▶ $Q(s, a) \leftarrow R(s, a) + \gamma \max_{a' \in A} Q(s', a')$
 - ▶ $s \leftarrow s'$
- ▶ Also known as “Dynamic Programming”, “Value Iteration”

35 / 70

AN EXAMPLE (1)



$$R(\text{HALL}, \text{To-CAVE}) = 0$$

$$Q(\text{CAVE}, a) = 0 \text{ for all actions } a$$

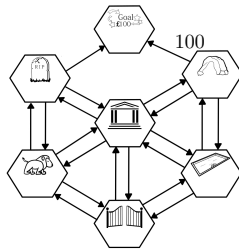
36 / 70

AN EXAMPLE (2)

Next suppose the agent, now in state CAVE, selects action $To - GOAL$

$R(CAVE, To-GOAL) = 100$, $Q(GOAL, a) = 0$ for all actions (there are no actions)

Hence $Q(CAVE, To-GOAL) = 100 + \gamma * 0 = 100$



37 / 70

AN EXAMPLE (3)

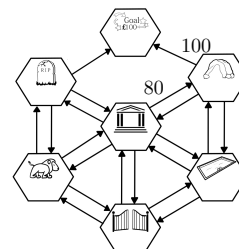
Let's start at hall again and select the same action To-CAVE

$R(HALL, To - CAVE) = 0$, $Q(CAVE, GOAL) = 100$

$Q(CAVE, a) = 0$ for all other actions a

Hence $\max_{a \in A} Q(CAVE, a) = 100$, if $\gamma = 0.8$,

$Q(HALL, To - CAVE) = 0 + \gamma * 100 = 80$



38 / 70

EXPLORATION / EXPLOITATION

- ▶ How do we best explore?
- ▶ Choose actions at random - but this can be very slow
- ▶ $\epsilon - greedy$ is the most common method
- ▶ Act ϵ -greedily
 - ▶ $\pi^\epsilon(s, a) = \begin{cases} a = \arg \max_{a \in A} Q(s, a) & \text{if } 1 - \epsilon + \epsilon/|A| \\ U_a & \text{otherwise} \end{cases}$
 - ▶ ϵ -greedy means acting greedily with probability $1 - \epsilon$, random otherwise
- ▶ When you are done, act greedily $\pi(s) = \arg \max_{a \in A} Q(s, a)$

39 / 70

ALGORITHMS FOR NON-DETERMINISTIC SETTINGS

- ▶ What can we do if the MDP is not deterministic?
- ▶ Q-learning
 - ▶ $Q(s, a) \leftarrow Q(s, a) + \eta \left[R(s, a) + \gamma \max_{a' \in A} Q(s', a') - Q(s, a) \right]$
- ▶ SARSA(0)
 - ▶ $Q(s, a) \leftarrow Q(s, a) + \eta [R(s, a) + \gamma Q(s', a') - Q(s, a)]$
- ▶ SARSA(1)/MC,
 - ▶ $Q(s, a) \leftarrow Q(s, a) + \eta [v_\tau - Q(s, a)]$
 - ▶ $v_\tau \leftarrow R(s, a) + \gamma R(s', a') + \dots \gamma^2 R(s'', a'') + \gamma^{\tau-1} R(s^\tau, a^\tau)$
- ▶ η is a small learning rate, e.g., $\eta = 0.001$

40 / 70

SARSA VS Q-LEARNING VS MC

- ▶ MC: updated using the whole chain
 - ▶ Possibly works better when the markov property is violated
- ▶ SARSA: update based on the next action you actually took
 - ▶ On Policy learning
- ▶ Q-Learning: update based on the best possible next action
 - ▶ Will learn optimal policy even if acting off-policy

41 / 70

MONTE CARLO CONTROL (1)

- ▶ Remember Q is just a mean/average
- ▶ MC (Naive Version)
 - ▶ Start at any state, initialise $Q_0(s, a)$ as you visit states/actions
 - ▶ Act ϵ -greedily
- ▶ Add all reward you have seen so far to $v_\tau^i = R(s', a') + \gamma R(s'', a'') + \gamma^2 R(s''', a''') + \gamma^{\tau-1} R(s^\tau, a^\tau)$ for episode i
- ▶ $Q_n(s, a) = E_{\pi^\epsilon}[v_\tau^i] = \frac{1}{n} \sum_{i=1}^n v_\tau^i$, where n is the times a state is visited

42 / 70

MONTE CARLO CONTROL (2)

- ϵ -greedy means acting greedily $1 - \epsilon$, random otherwise
- Better to calculate mean incrementally

$$Q_n(s, a) = E_{\pi_n}[v_\tau^i]$$

$$Q_n(s, a) = \frac{1}{n} \sum_{i=1}^n v_\tau^i$$

$$Q_n(s, a) = \frac{1}{n} (v_\tau^1 + v_\tau^2 \dots v_\tau^{n-1} + v_\tau^n)$$

$$Q_n(s, a) = \frac{1}{n} \left(\sum_{i=1}^{n-1} v_\tau^i + v_\tau^n \right)$$

43 / 70

MONTE CARLO CONTROL (3)

by definition

$$Q_{n-1}(s, a) = \frac{1}{n-1} \sum_{i=1}^{n-1} v_\tau^i \implies (n-1) Q_{n-1}(s, a) = \sum_{i=1}^{n-1} v_\tau^i$$

$$Q_n(s, a) = \frac{1}{n} ((n-1) Q_{n-1}(s, a) + v_\tau^n)$$

$$Q_n(s, a) = \frac{1}{n} (Q_{n-1}(s, a) n - Q_{n-1}(s, a) + v_\tau^n)$$

$$Q_n(s, a) = \frac{Q_{n-1}(s, a) n}{n} + \frac{-Q_{n-1}(s, a) + v_\tau^n}{n}$$

$$Q_n(s, a) = Q_{n-1}(s, a) + \frac{\overbrace{v_\tau^n - Q_{n-1}(s, a)}^{\text{MC-Error}}}{n}$$

44 / 70

MONTE CARLO CONTROL (4)

- But π^n changes continuously, so the distribution of rewards is non-stationary

$$Q_n(s, a) = Q_{n-1}(s, a) + \frac{1}{n} [v_\tau^n - Q_{n-1}(s, a)] \rightarrow \text{Bandit case}$$

$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [v_\tau^n - Q_{n-1}(s, a)] \rightarrow \text{Full MDP case}$$

- A Bandit can be seen as MDP with a chain of length one (i.e. s) - η is a learning rate (e.g., 0.001)

45 / 70

MONTE CARLO CONTROL (5)

- Start at any state, initialise $Q_0(s, a)$ as you visit states/actions
- Act ϵ -greedily
- Wait until episode ends, i.e. a terminal state is hit - ϵ set to some low value, e.g., 0.1
- Add all reward you have seen so far to $v_\tau^i = R(s, a) + \gamma R(s', a') + \dots \gamma^2 R(s'', a'') + \gamma^{\tau-1} R(s^\tau, a^\tau)$ for episode i
- $Q_n(s, a) = Q_{n-1}(s, a) + \eta [v_\tau^n - Q_{n-1}(s, a)]$

46 / 70

FROM MONTE CARLO CONTROL TO SARSA AND Q-LEARNING

- With MC we update using the rewards from the whole chain
- Can we update incrementally?

$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [v_\tau^n - Q_{n-1}(s, a)]$$

$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [R(s, a) + \gamma R(s', a') + \dots \gamma^2 R(s'', a'') + \gamma^{\tau-1} R(s^\tau, a^\tau) - Q_{n-1}(s, a)]$$

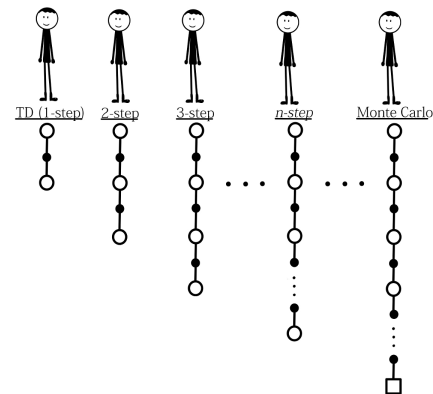
$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [R(s, a) + \gamma (R(s', a') + \dots \gamma R(s'', a'') + \gamma^{\tau-2} R(s^\tau, a^\tau)) - Q_{n-1}(s, a)]$$

$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [R(s, a) + \gamma (v_\tau^n(s', a') - Q_{n-1}(s, a))]$$

$$Q_n(s, a) = Q_{n-1}(s, a) + \eta [R(s, a) + \gamma Q_{n-1}(s', a') - Q_{n-1}(s, a)]$$

47 / 70

N-STEP RETURNS



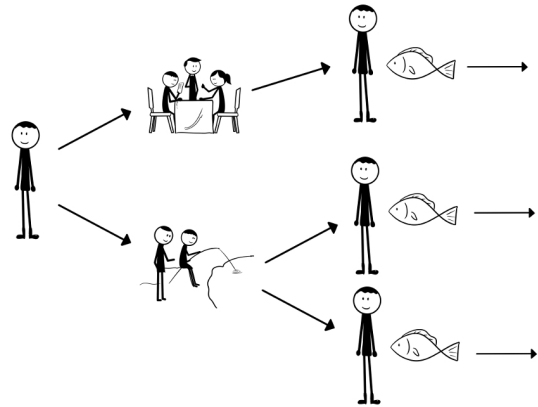
48 / 70

LET'S GO OVER THE TOON EXAMPLE, WITHOUT A MODEL

- ϵ – greedy, with $\epsilon = 0.1$

49 / 70

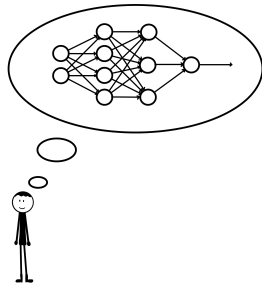
MODEL FREE TOON



50 / 70

FUNCTION APPROXIMATION (1)

- There is usually some link between states
- We can train function approximators incrementally to model $Q(s, a)$
- We now have $Q(s, a; \theta)$, where θ are the parameters



51 / 70

FUNCTION APPROXIMATION (2)

- What are the links in states in Toon?
- Can we write down the Q-values in a more compact way?
 - Let's devise a method to do this
- Examples include linear function approximators, neural networks, n-tuple networks
- Not easy to do, few convergence guarantees
 - But with some effort, this works pretty well

52 / 70

POLICY WITH FEATURES

- What if after catching fish there was another action to choose from ("how many should I eat?")

Policy	Policy Value	Q-Values
$\pi(\text{Start, Go-Fishing})$?	?
$\pi(\text{Start, Go-to-Restaurant})$?	?
$\pi(\text{Restaurant, Eat-}\phi\text{-fish})$	1	ϕ

53 / 70

WHAT DO WE ACTUALLY LEARN?

- X are our features
- Targets are
 - Q-learning
 - $y = R(s, a) + \gamma \max_{a' \in A} Q(s', a')$
 - SARSA(0)
 - $y = R(s, a) + \gamma Q(s', a')$
 - SARSA(1)/MC,
 - $y \leftarrow v_\tau$
 - $v_\tau \leftarrow R(s, a) + \gamma R(s', a') + \dots \gamma^2 R(s'', a'') + \gamma^{\tau-1} R(s^\tau, a^\tau)$
- N-Step versions
 - Same as MC version, but stop prematurely and take a SARSA/Q-learning target

54 / 70

INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 275 558 302">WHAT CAN BE USED AS FEATURES?</h2> <ul data-bbox="144 422 753 596" style="list-style-type: none"> ▶ Anything (text, sound chunks, images) ▶ For text see here: <ul style="list-style-type: none"> ▶ https://github.com/facebookresearch/CommAI-env ▶ You often don't need to start from scratch, for text you have <i>word2vec</i> ▶ Different Neural Network architectures <p data-bbox="753 758 797 772">55 / 70</p>	<h2 data-bbox="831 275 1300 338">NEURAL NETWORKS AND FUNCTION APPROXIMATION</h2> <ul data-bbox="876 443 1500 611" style="list-style-type: none"> ▶ Most common modern function approximation scheme is neural networks ▶ Can approximate almost any function ▶ We had a series of recent advances <ul style="list-style-type: none"> ▶ Go (10^{170} states) ▶ Atari (grayscale, 110 x 84 resolution) <p data-bbox="1484 758 1528 772">56 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 827 250 854">PLATFORMS</h2> <ul data-bbox="144 926 721 1205" style="list-style-type: none"> ▶ Tools <ul style="list-style-type: none"> ▶ Keras (neural networks) ▶ Tensorflow (neural networks, but closer to the machine) ▶ <code>goo.gl/YGWSbL</code> ▶ Open AI gym ▶ There is a phenomenal lack of windows support! ▶ Let's look at open AI gym ▶ A lot of modern work is a combination of RL with neural networks ▶ We have good libraries now <p data-bbox="753 1304 797 1318">57 / 70</p>	<h2 data-bbox="831 827 1219 854">MORE ON NEURAL NETWORKS</h2> <ul data-bbox="876 953 1484 1163" style="list-style-type: none"> ▶ A function approximator loosely based on the brain ▶ Global function approximator ▶ Catastrophic forgetting... ▶ Multiple ways of breaking correlations <ul style="list-style-type: none"> ▶ Experience replay, asynchronous games ▶ Again, think of Neural Networks as a mechanism for storing Q-Values <p data-bbox="1484 1304 1528 1318">58 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 1379 444 1407">WHAT ARE WE LEARNING?</h2>	<h2 data-bbox="831 1379 1289 1407">NEURAL NETWORK ARCHITECTURE</h2> <ul data-bbox="876 1514 1354 1703" style="list-style-type: none"> ▶ There are certain choices that need to be made ▶ Number of layers ▶ Type of layers ▶ Learning algorithms ▶ Regularisation methods ▶ Many different ways of building those networks ▶ Let's look at some code <p data-bbox="1484 1850 1528 1864">60 / 70</p>
59 / 70	

INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="94 275 363 302">INTUITION BUILDING</h2> <ul data-bbox="142 449 690 554" style="list-style-type: none"> ▶ Choose a game ▶ Choose a character in the game ▶ Chose the features that represent the character's state ▶ Choose the neural network to use <p data-bbox="751 753 797 770">61 / 70</p>	<h2 data-bbox="828 275 1127 302">SINGLE PLAYER GAMES</h2> <ul data-bbox="873 407 1498 611" style="list-style-type: none"> ▶ Everything we have seen is based on single player environments <ul data-bbox="919 443 1498 491" style="list-style-type: none"> ▶ But from NPC perspective there is no such thing as single player ▶ The actual player is your opponent! ▶ Domain of multiple agents interacting is <i>Game Theory</i> (or multi-agent learning) ▶ Environment adapts back at you ▶ Needs more tricks to get things to perform sensibly <p data-bbox="1482 753 1528 770">62 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="94 827 789 854">RELATIONSHIP TO THE REST OF MACHINE LEARNING</h2> <ul data-bbox="142 917 768 1220" style="list-style-type: none"> ▶ How can one learn a model of the world? <ul data-bbox="188 953 727 1073" style="list-style-type: none"> ▶ Possibly by breaking it down into smaller, abstract chunks <ul data-bbox="233 984 443 1008" style="list-style-type: none"> ▶ Unsupervised Learning ▶ ... and learning what effects ones actions have the environment <ul data-bbox="233 1050 422 1073" style="list-style-type: none"> ▶ Supervised Learning ▶ RL weaves all fields of Machine Learning (and possibly Artificial Intelligence) into one coherent whole ▶ The purpose of all learning is action! <ul data-bbox="188 1176 755 1220" style="list-style-type: none"> ▶ You need to be able to recognise faces so you can create state ▶ ... and act on it <p data-bbox="751 1304 797 1320">63 / 70</p>	<h2 data-bbox="828 827 1279 854">CAUSALITY (A VERY BRIEF INTRO)</h2> <ul data-bbox="873 1016 1328 1068" style="list-style-type: none"> ▶ We often colloquially say “A is caused by B” ▶ Can you discuss the meaning of this? <p data-bbox="1482 1304 1528 1320">64 / 70</p>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="94 1379 352 1407">COUNTERFACTUALS</h2> <ul data-bbox="142 1526 644 1682" style="list-style-type: none"> ▶ If I take action a I land on state s ▶ What if I don't take action a? ▶ “Experimenter forced you to pick up smoking” vs ▶ “Experimenter observed that you smoked” ▶ Will you get lung disease? ▶ The experimenter takes the actions vs observes <p data-bbox="751 1850 797 1866">65 / 70</p>	<h2 data-bbox="828 1379 1089 1407">WHAT IS THE LINK?</h2> <ul data-bbox="873 1520 1446 1690" style="list-style-type: none"> ▶ Off-policy evaluation learning ▶ Let's see an example <ul data-bbox="919 1581 1325 1627" style="list-style-type: none"> ▶ Features are colour of hair, height, smoking ▶ Reward is 0 (lung disease), 1 (healthy) ▶ This would have been supervised learning if we knew the policy! <p data-bbox="1482 1850 1528 1866">66 / 70</p>

INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 275 261 302">CONCLUSION</h2> <ul data-bbox="142 449 756 552" style="list-style-type: none"> ▶ RL is a massive topic ▶ We have shown the tip of iceberg ▶ Rabbit hole goes <i>deep</i> - both on the application level and the theory level <div data-bbox="753 756 797 772">67 / 70</div>	<h2 data-bbox="829 275 1094 302">FURTHER STUDY (1)</h2> <ul data-bbox="873 378 1500 659" style="list-style-type: none"> ▶ Tom Mitchell, Chapter 13 ▶ David Silver's UCL Course: <ul data-bbox="898 430 1500 533" style="list-style-type: none"> ▶ http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html <ul data-bbox="922 468 1500 533" style="list-style-type: none"> ▶ Some ideas in these lecture notes taken from there ▶ Probably the best set of notes there is on the subject ▶ Online at http://www.machinelearningtalks.com/tag/rl-course/ ▶ Reinforcement Learning, by Richard S. Sutton and Andrew G. Barto <ul data-bbox="922 613 1265 659" style="list-style-type: none"> ▶ Classic book ▶ Excellent treatment of most subjects <div data-bbox="1482 756 1526 772">68 / 70</div>
INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING	INTRODUCTION & MOTIVATION · MARKOV DECISION PROCESS (MDPs) · PLANNING · MODEL FREE REINFORCEMENT LEARNING
<h2 data-bbox="99 827 363 854">FURTHER STUDY (2)</h2> <ul data-bbox="142 919 768 1220" style="list-style-type: none"> ▶ Artificial Intelligence: A Modern Approach by Stuart J. Russell and Peter Norvig <ul data-bbox="188 980 493 1026" style="list-style-type: none"> ▶ The Introductory A.I. Textbook ▶ Chapters 16 and 21 ▶ Algorithms for Reinforcement Learning by Csaba Szepesvari <ul data-bbox="188 1077 768 1123" style="list-style-type: none"> ▶ Very “Mathematical”, but a good resource that provides a very unified view of the field ▶ Reinforcement Learning: State-Of-The-Art by Marco Wiering (Editor), Martijn Van Otterlo (Editor) <ul data-bbox="188 1197 342 1220" style="list-style-type: none"> ▶ Edited Volume <div data-bbox="753 1304 797 1320">69 / 70</div>	<h2 data-bbox="829 827 1127 854">SOME MODERN PAPERS</h2> <ul data-bbox="873 974 1487 1136" style="list-style-type: none"> ▶ Asynchronous Methods for Deep Reinforcement Learning https://arxiv.org/pdf/1602.01783v2.pdf ▶ A Survey of Monte Carlo Tree Search Methods http://www.cameronius.com/cv/mcts-survey-master.pdf ▶ Deep Exploration via Bootstrapped DQN https://arxiv.org/abs/1602.04621 <div data-bbox="1482 1304 1526 1320">70 / 70</div>