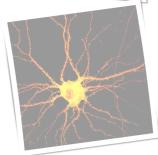
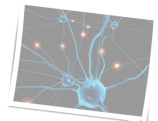
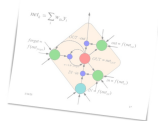
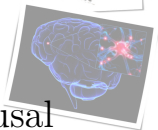


An introduction to Reinforcement Learning (with an intro to neural networks and causal reasoning)

Spyros Samothrakis
Research Fellow, IADS
Univerisity of Essex

November 14, 2016



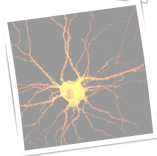
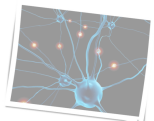
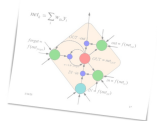
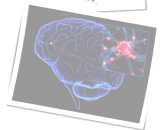
Introduction & Motivation

Markov Decision Process (MDPs)

Planning

Model Free Reinforcement Learning

Causality



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

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

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

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

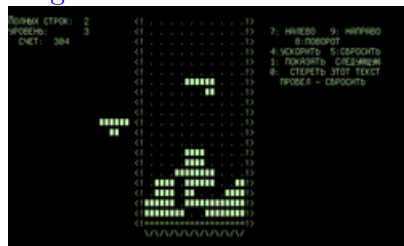
THE MARKOV PROPERTY AND STATES

- ▶ 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



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)

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

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”?

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?

LONG TERM THINKING

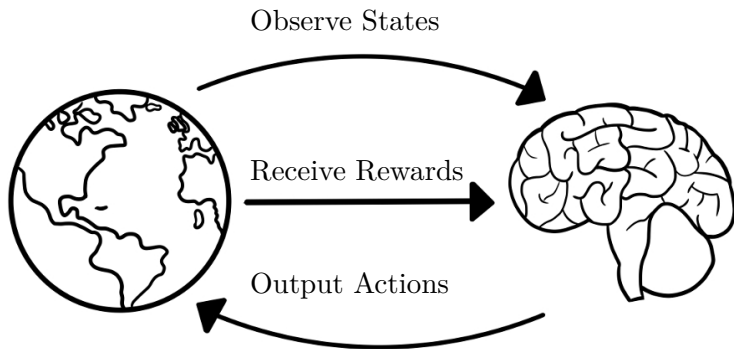
- ▶ It might be better to delay satisfaction
- ▶ Immediate reward is not always the maximum reward
- ▶ In some settings there are no immediate rewards at all (e.g., most solitaire games)
- ▶ MDPs and RL capture this
- ▶ “Not going out tonight, study”
- ▶ Long term investment

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

THE FULL LOOP

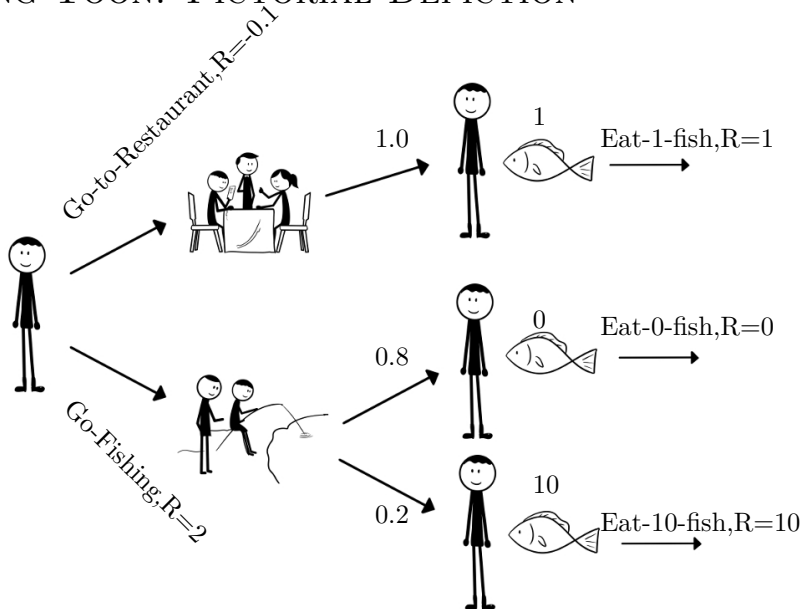
- ▶ See how the universe described by the MDP defines actions, not just states and transitions
- ▶ An agent needs to act upon what it perceives
- ▶ Notice the lack of body - “brain in a vat”. Body is assumed to be part of the world.



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

FISHING TOON: PICTORIAL DEPICTION

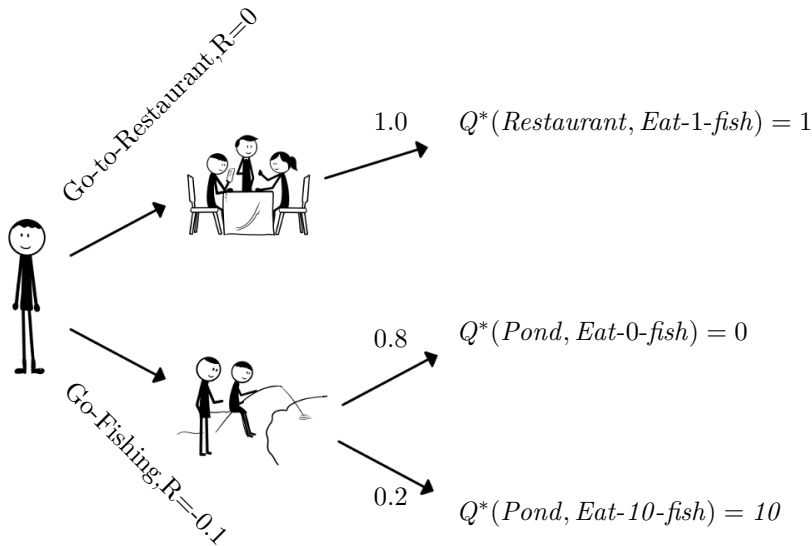


SUM OF EXPECTED REWARDS

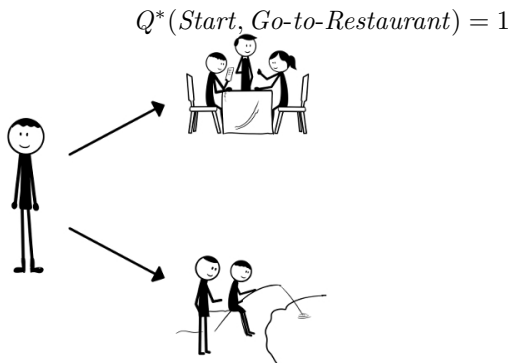
- ▶ Our toon has to choose between two different actions
- ▶ **Go-To-Restaurant** or **Go-Fishing**
- ▶ We assume that toon is interested in maximising *the expected sum* of happiness/reward
- ▶ Let's first see what happens if we start with a random policy

| 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 | |

RANDOM POLICY (1)



RANDOM POLICY (2)



$$Q^*(Start, Go-Fishing) = 0.2 * 10 + 0. * 8 * 0.0 - 0.1 = 1.9$$

TABLE

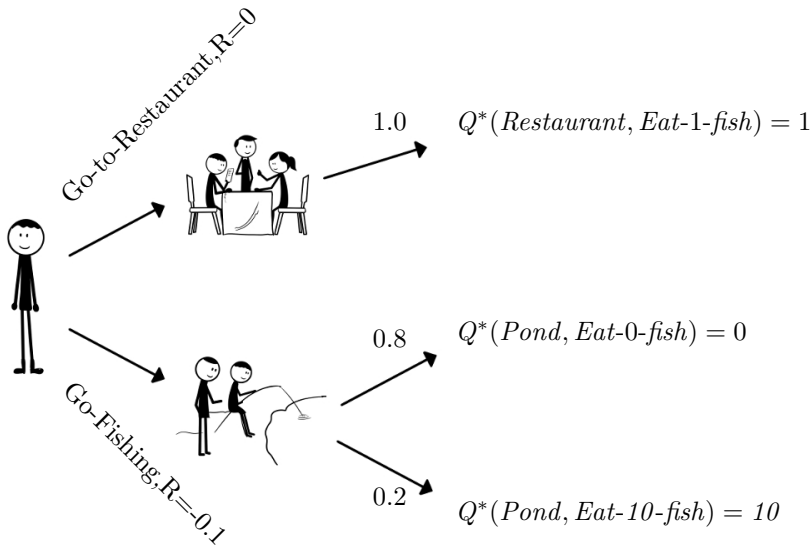
| Policy | Policy Value | Q-Values |
|--------------------------------|--------------|----------|
| $\pi(Start, Go-Fishing)$ | 0.5 | 1 |
| $\pi(Start, Go-to-Restaurant)$ | 0.5 | 1.9 |
| $\pi(Restaurant, Eat-1-fish)$ | 1 | 1 |
| $\pi(Pond, Eat-0-fish)$ | 1 | 0 |
| $\pi(Pond, Eat-10-fish)$ | 1 | 10 |

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

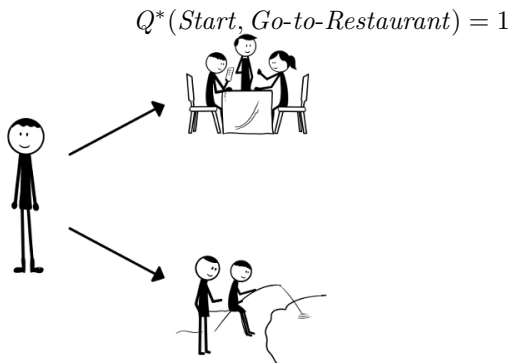
WHAT IF WE ARE ASKED TO FIND OUT THE OPTIMAL POLICY?

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

REASONING BACKWARDS (1)



REASONING BACKWARDS (2)



$$Q^*(Start, Go-Fishing) = 0.2 * 10 + 0. * 8 * 0.0 - 0.1 = 1.9$$

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 $V^*(Start) = \max\{1, 1.9\} = 1.9$

CORRECT ACTION

- ▶ Toon should go Go-Fishing
- ▶ **Would you do the same?**
- ▶ **Would a pessimist toon do the same?**
- ▶ We just went through the following equation:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) \max_{a' \in A} Q^*(s', a')$$

- ▶ Looks intimidating - but it's really simple
- ▶ Let's have a look at another example
 - ▶ How about toon goes to the restaurant after failing to fish?
 - ▶ How would that change the reward structure?

AGENT GOALS

- ▶ 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?
 - ▶ Sociopathic, egotistic, greed-is-good Gordon Gekko (*Wall Street*, 1987)
 - ▶ No concept of “externalities” - agents might wreak havoc for marginal reward gains
 - ▶ Same applies to all “compulsive agents” - think Chess

SEARCHING FOR A GOOD POLICY

- ▶ One can possibly search through all combinations of policies until she finds the best
- ▶ Slow, does not work in larger MDPs
- ▶ Exploration/Exploitation dilemma
 - ▶ How much time/effort should be spend exploring for solutions?
 - ▶ How much time should be spend exploiting good solutions?

PLANNING

- ▶ 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”

BELLMAN EXPECTATION EQUATIONS / BELLMAN BACKUPS

- ▶ 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

LINK BETWEEN V^π AND Q^π

- ▶ 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!**

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
- ▶ $V^*(s) = \max_{a \in A} R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) V^*(s')$
- ▶ $Q^*(s, a) = R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) \max_{a' \in A} Q^*(s', a')$
- ▶ $\pi^*(s, a) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} Q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$

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*

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)...

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

$$\text{▶ } T(s'|s, a_i) = \begin{cases} 1 & \text{if } a_i = a_c \\ 0 & \text{otherwise} \end{cases}$$

$$\text{▶ } V^*(s) = \max_{a \in A} [R(s, a) + \gamma V^*(s')]$$

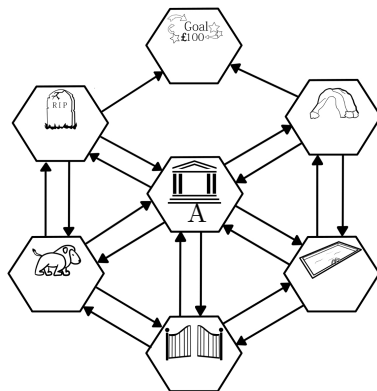
$$\text{▶ } Q^*(s, a) = R(s, a) + \gamma \max_{a' \in A} Q(s', a')$$

- ▶ 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 $Q(s, a)$

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”

AN EXAMPLE (1)



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

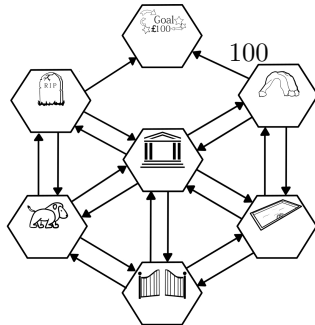
$$Q(\text{CAVE}, a) = 0 \text{ for all actions } a$$

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



AN EXAMPLE (3)

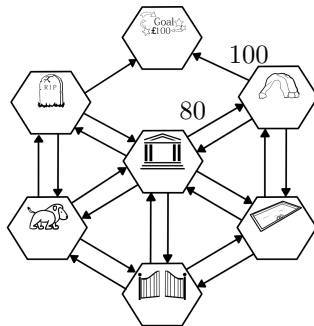
Let's start at hall again and select the same action To-CAVE

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

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

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

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



EXPLORATION / EXPLOITATION

- ▶ How do we best explore?
- ▶ Choose actions at random - but this can be very slow
- ▶ ϵ - *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 $\pi(s) = \arg \max_{a \in A} Q(s, a)$

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., $\eta = 0.001$

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

MONTÉ 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

MONTÉ CARLO CONTROL (2)

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

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

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

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

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

MONTE CARLO CONTROL (3)

by definition

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

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

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

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

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

MONTÉ CARLO CONTROL (4)

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

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

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

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

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

FROM MONTE CARLO CONTROL TO SARSA AND Q-LEARNING

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

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

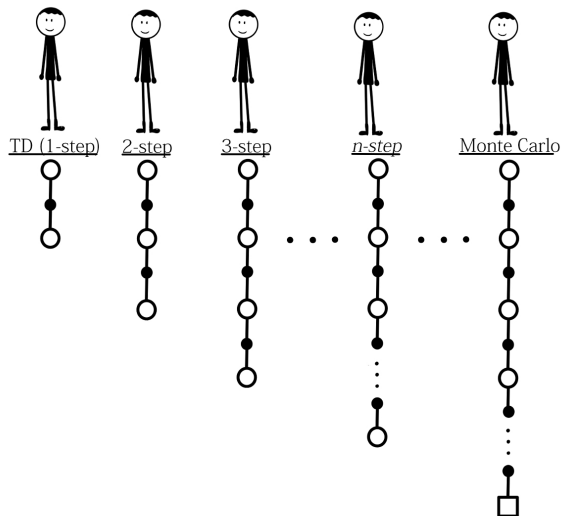
$$Q_n(s, a) = Q_{n-1}(s, a) + \eta \left[R(s, a) + \gamma R(s', a') + \dots \gamma^2 R(s'', a'') + \gamma^{\tau-1} R(s^\tau, a^\tau) - Q_{n-1}(s, a) \right]$$

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

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

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

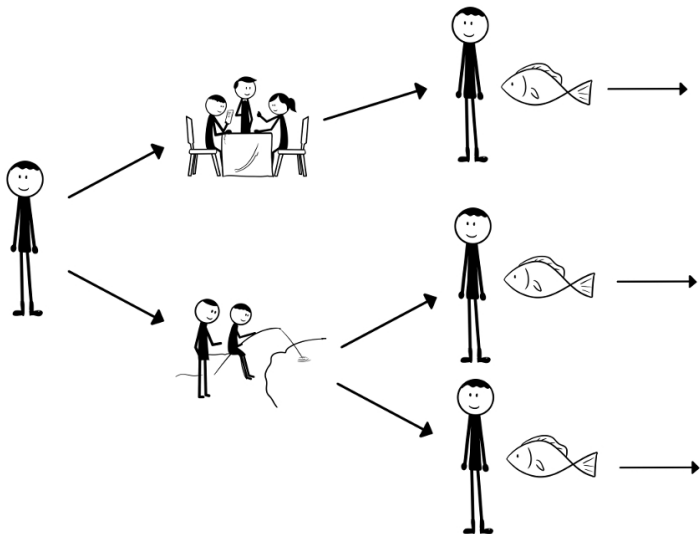
N-STEP RETURNS



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

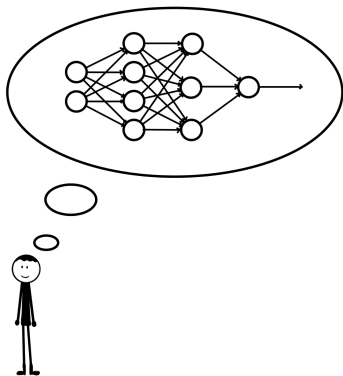
- ▶ ϵ - *greedy*, with $\epsilon = 0.1$

MODEL FREE TOON



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



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 tree 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

POLICY WITH FEATURES

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

DO WE HAVE TO LEARN Q-VALUES?

► ?

NEURAL NETWORKS AND FUNCTION APPROXIMATION

- ▶ Most common modern function approximation scheme is neural networks
- ▶ Can approximate almost any function
- ▶ We had a series of recent advances
 - ▶ Go (10^{170} states)
 - ▶ Atari (10^{10^7} states)

PLATFORMS

- ▶ Tools
- ▶ Keras (neural networks)
- ▶ Tensorflow (neural networks, but closer to the machine)
- ▶ goo.gl/YGWSbL
- ▶ Open AI gym
- ▶ 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

MORE ON NEURAL NETWORKS

- ▶ A function approximator loosely based on the brain
- ▶ Global function approximator
- ▶ Catastrophic forgetting....
- ▶ Multiple ways of breaking correlations
 - ▶ Experience replay, asynchronous games
- ▶ Again, think of Neural Networks as a mechanism for storing Q-Values

NEURAL NETWORK ARCHITECTURE

- ▶ 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

INTUITION BUILDING

- ▶ Choose a game
- ▶ Choose a character in the game
- ▶ Choose the features that represent the character's state
- ▶ Choose the neural network to use

SINGLE PLAYER GAMES

- ▶ Everything we have seen is based on single player environments
 - ▶ But from NPC perspective there is no such thing as single player
 - ▶ The actual player is your opponent!
- ▶ Domain of multiple agents interacting is *Game Theory* (or multi-agent learning)
- ▶ Environment adapts back at you
- ▶ Needs more tricks to get things to perform sensibly

RELATIONSHIP TO THE REST OF MACHINE LEARNING

- ▶ How can one learn a model of the world?
 - ▶ Possibly by breaking it down into smaller, abstract chunks
 - ▶ Unsupervised Learning
 - ▶ ... and learning what effects ones actions have the environment
 - ▶ Supervised Learning
- ▶ RL weaves all fields of Machine Learning (and possibly Artificial Intelligence) into one coherent whole
- ▶ The purpose of all learning is action!
 - ▶ You need to be able to recognise faces so you can create state
 - ▶ ... and act on it

CAUSALITY (A VERY BRIEF INTRO)

- ▶ We often colloquially say “A is caused by B”
- ▶ Can you discuss the meaning of this?

COUNTERFACTUALS

- ▶ 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

WHAT IS THE LINK?

- ▶ Off-policy evaluation learning
- ▶ Let's see an example
 - ▶ Features are colour of hair, height, smoking
 - ▶ Reward is -1000 (lung disease), 1 (healthy)
- ▶ This would have been supervised learning if we knew the policy!

CONCLUSION

- ▶ RL is a massive topic
- ▶ We have shown the tip of iceberg
- ▶ Rabbit hole goes *deep* - both on the application level and the theory level

FURTHER STUDY (1)

- ▶ **Tom Mitchell, Chapter 13**
- ▶ David Silver's UCL Course:
<http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html>
 - ▶ 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
 - ▶ Classic book
 - ▶ Excellent treatment of most subjects

FURTHER STUDY (2)

- ▶ Artificial Intelligence: A Modern Approach by Stuart J. Russell and Peter Norvig
 - ▶ The Introductory A.I. Textbook
 - ▶ Chapters 16 and 21
- ▶ Algorithms for Reinforcement Learning by Csaba Szepesvari
 - ▶ 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)
 - ▶ Edited Volume