

5a: Reinforcement Learning

Week 5: Overview

The topic for this week is Reinforcement Learning (RL). We will see how reinforcement learning tasks can be formally defined, compare different models of optimality, and discuss the need for exploration and the problem of delayed rewards. We will present a number of RL algorithms including Temporal Difference learning (TD-learning), Q-learning, policy gradients, actor-critic. We will see how RL can be applied to games like backgammon, and how deep Q-learning or A3C can learn to play Atari games from raw pixels.

Weekly learning outcomes

By the end of this week, you will be able to:

- explain the difference between supervised learning, reinforcement learning, and unsupervised learning
 - describe the formal definition of a reinforcement learning task
 - identify different models of optimality
 - explain the need for exploration in reinforcement learning
 - describe reinforcement learning algorithms, including TD-learning, Q-learning, policy gradients, and actor-critic
 - apply Q-learning to simple tasks.
-

Reinforcement Learning

We have previously discussed supervised learning, where pairs of inputs and target outputs are provided and the system must learn to predict the correct output for each input.

There are many situations where we instead want to train a system to perform certain actions in an environment, in order to maximise a reward function. These situations include, for example, playing a video game, allocating mobile phone channels or other resources dynamically, driving a car or flying a helicopter.

Supervised learning can sometimes be used for this purpose, if we construct a training set of situation-action pairs (for example, a database of game positions and the move chosen by a human expert, or sensor readings from a motor car and the steering direction chosen by a human driver). This process is sometimes called Behavioral Cloning.

However, it is often better if the system can learn by purely by self-play, without the need for training data from a human expert.

Reinforcement Learning Framework

Reinforcement Learning (RL) can be formalised in terms of an Agent interacting with its Environment.

The **Environment** includes a set S of states and a set A of actions. At each time step t , the agent is in some state s_t . It must choose an action a_t , whereupon it goes into state $s_{t+1} = \delta(s_t, a_t)$ and receives reward $r_t = R(s_t, a_t)$

The **Agent** chooses its actions according to some **policy** $\pi : S \rightarrow A$.

The aim of Reinforcement Learning is to find an **optimal** policy π^* which maximises the cumulative reward.

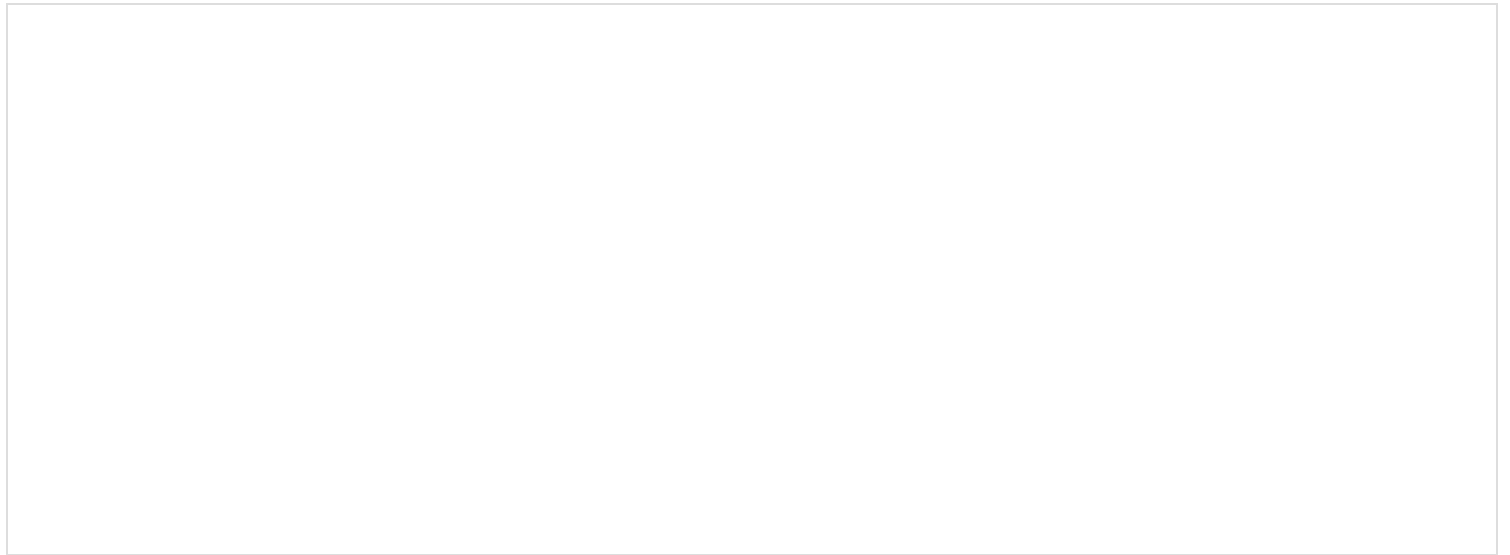
In some cases, the Environment may be probabilistic or **stochastic** — meaning that the transitions and/or rewards are not fully determined, but instead occur randomly according to some probability distribution:

$$\delta(s_{t+1} = s \mid s_t, a_t), \quad R(r_t = r \mid s_t, a_t)$$

The Agent can also employ a probabilistic policy, with its actions chosen according to a probability distribution:

$$\pi(a_t = a \mid s_t)$$

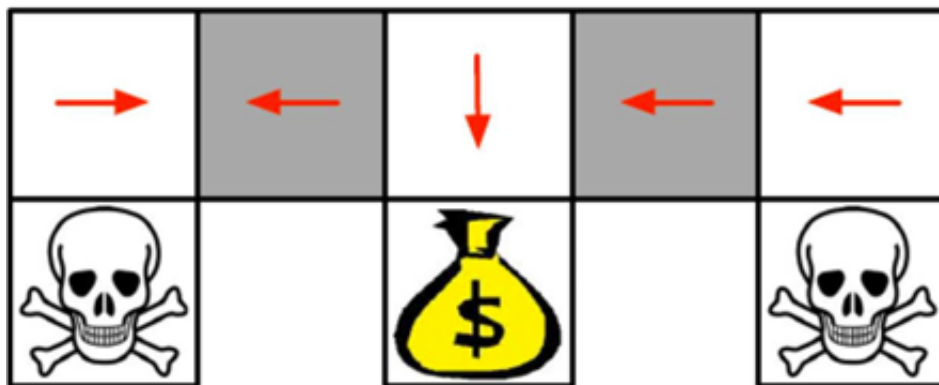
Optional video



Probabilistic Policies

The main benefit of a probabilistic policy is that it forces the agent to explore its environment more comprehensively while it is learning.

However, probabilistic policies may have additional benefits in other contexts, such as partially observable or multi-player environments.



Consider, for example, a situation in which the Agent is able to observe that it is in one of the grey squares, but is not able to determine which grey square it is in. In this case, the policy of moving left or right with equal probability from the grey square will perform better than any deterministic policy.

In two-player games like rock-paper-scissors, a random strategy is also required in order to make agent choices unpredictable to the opponent.

Optional video

Models of Optimality

What exactly do we mean when we say that an RL Agent should choose a policy which maximises its future rewards?

The old saying "a fast nickel is worth a slow dime" reminds us that people often prefer to receive a smaller reward sooner instead of a larger reward later. Economists use surveys to gauge people's preferences in this regard: "Would you prefer ten dollars today, or 15 dollars next week? How about 50 dollars today compared to 100 dollars next year?" Responses to these surveys can be used to estimate a number γ between 0 and 1 that is chosen such that one dollar in the current timestep is considered equally desirable with γ dollars in the next timestep. In Reinforcement Learning, this number γ is called the **discount factor** and is used to define the infinite discounted reward, which can be compared with average reward and finite horizon reward.

$$\text{Finite horizon reward: } \sum_{i=0}^{h-1} r_{t+i}$$

$$\text{Infinite discounted reward: } \sum_{i=0}^{\infty} \gamma^i r_{t+i}, \quad 0 \leq \gamma < 1$$

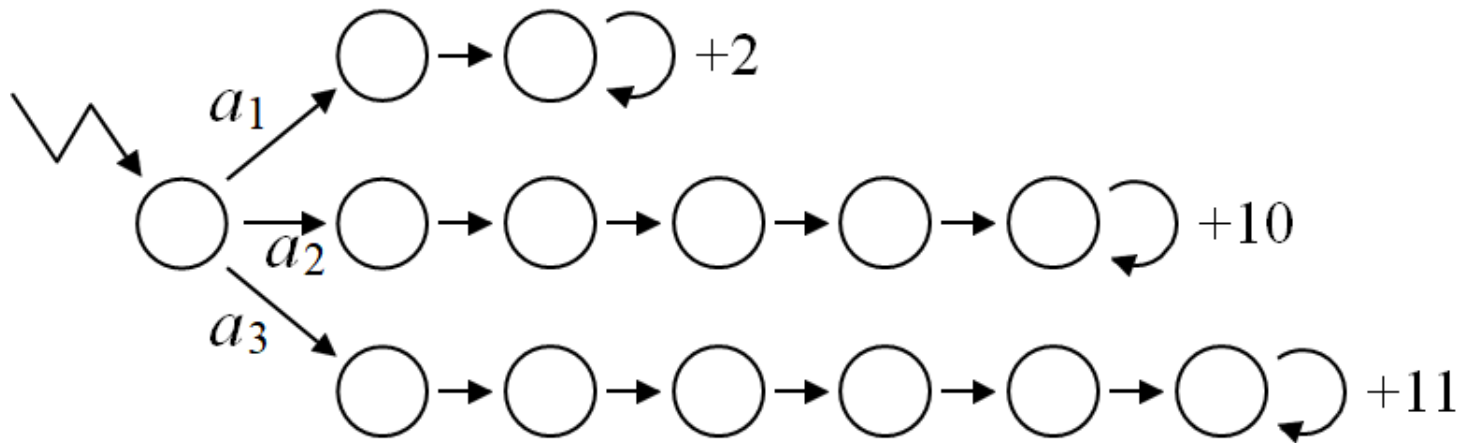
$$\text{Average reward: } \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{i=0}^{h-1} r_{t+i}$$

We normally try to maximise the infinite discounted reward, because it is easier to analyse theoretically and also works well in practice.

The finite horizon reward is easy to compute, but may lead to bad decisions by failing to take into account rewards which are just over the horizon.

Average reward is hard to deal with because we cannot sensibly choose between a small reward soon and a large reward very far in the future — for example, 100 dollars today compared with a million dollars, 50 years from now.

Comparing Models of Optimality



This environment illustrates how the choice of action may depend on the model of optimality. An agent trying to maximise the finite horizon reward with $h=4$ would prefer action a_1 because it is the only action which will gain any reward within the first four time steps. An agent maximising average reward will prefer action a_3 because, after the first few time steps, it will be getting a reward of +11 every timestep which is larger than +10 for an agent who chose action a_2 . An agent maximising infinite discounted reward with $\gamma = 0.9$ will prefer action a_2 . One way to see this is as follows: Consider Agent A_2 choosing action a_2 , compared to Agent A_3 choosing action a_3 . The reward of 11 received by A_3 in timestep 7 would be equivalent to 9.9 in timestep 6 and therefore slightly less than the \$10 received by A_2 in timestep 6. By the same logic, the \$10 received by A_2 in timestep n is always slightly more valuable than the \$11 received by A_3 in timestep $n+1$.

Optional video



Value Function

Every policy π determines a Value Function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ where $V^\pi(s)$ is the average discounted reward received by an agent who begins in state s and chooses its actions according to policy π .

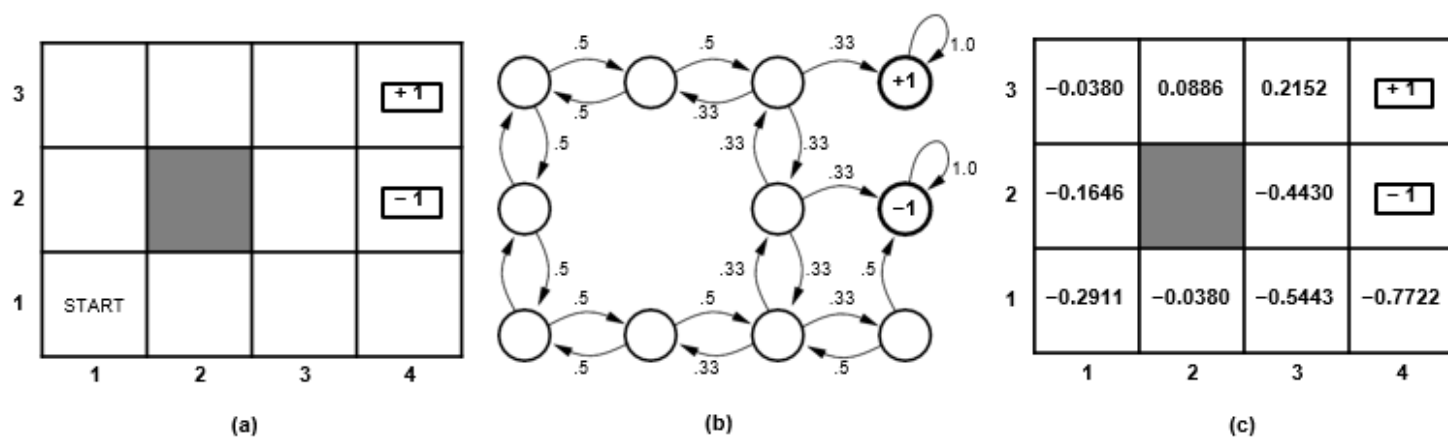
For small environments, we can compute the value function using simultaneous equations. For larger environments, we need a way of learning it, through experience.

Transcript

If $\gamma = 0.9$, the value of the last node in the middle row of the above diagram is $\frac{10}{1-\gamma} = 100$. The value of the last node in the lowest row is $\frac{11}{1-\gamma} = 110$. The value of the second last node in the lower row is 99, which can be obtained from the value of the last node either by subtracting 11 or multiplying by γ .

Value Function Example

This image shows the value function V_π in a grid world environment, where π is the policy of choosing between available actions uniformly randomly.



Exploration and Delayed Reinforcement

I was born to try ...
But you've got to make choices
Be wrong or right
Sometimes you've got to sacrifice the things you like.
— Delta Goodrem

Multi-Armed Bandit Problem



The special case of a stochastic reinforcement learning environment with only one state is called the **Multi-Armed Bandit Problem**, because it is like being in a room with several (friendly) slot machines (also called "one-armed bandits") for a limited time, and trying to collect as much money as possible. We assume that each slot machine provides rewards chosen from its own (stationary) distribution, and that in each time step we are able to pull the lever on only one machine.

Exploration / Exploitation Tradeoff

After pulling the levers a few times and observing the rewards received, it makes sense that most of the time we should choose the lever (action) which we think will give the highest expected reward, based on previous observations.

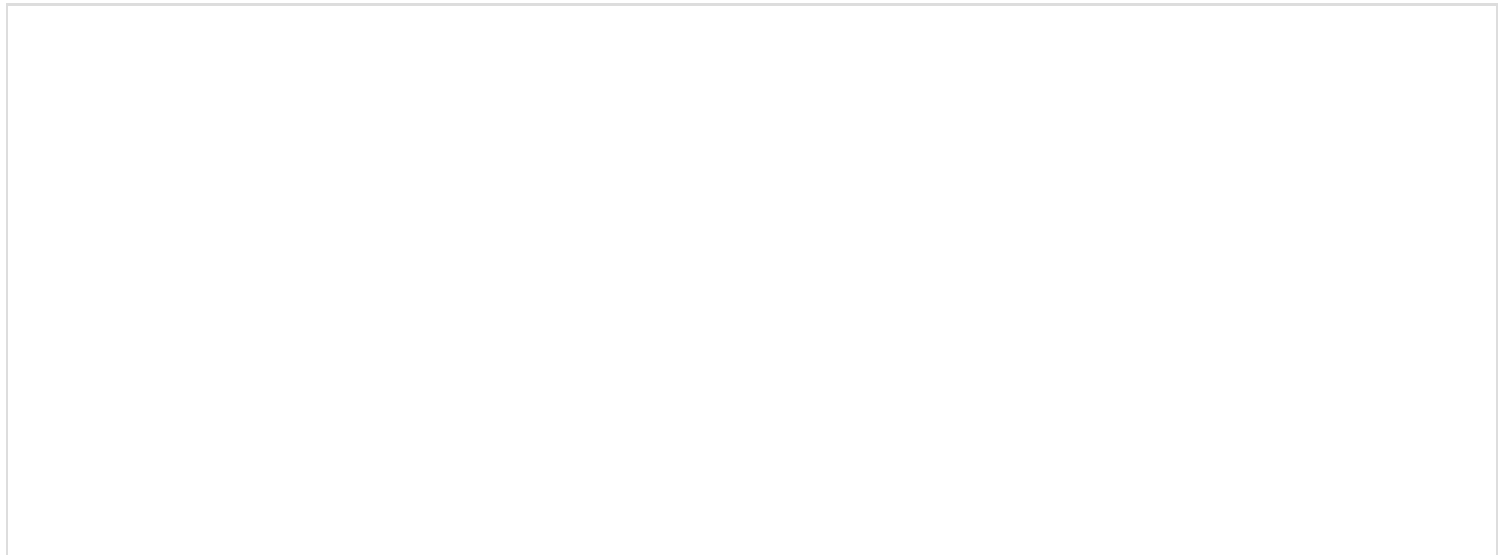
However, in order to ensure convergence to the optimal strategy, we must occasionally choose something different from our preferred action. Perhaps the simplest way to achieve this is known as an **epsilon-greedy strategy**, where with probability $1 - \epsilon$ we choose what we think is the best action, and with probability ϵ (typically, 5%) we choose a random action. More sophisticated strategies also exist such as Thompson Sampling, Upper Confidence Bound (UCB) algorithms, or choosing from a Softmax (Boltzmann) distribution based on the average reward so far observed for

each action.

$$P(a) = \frac{e^{R(a)/T}}{\sum_{b \in A} e^{R(b)/T}}$$

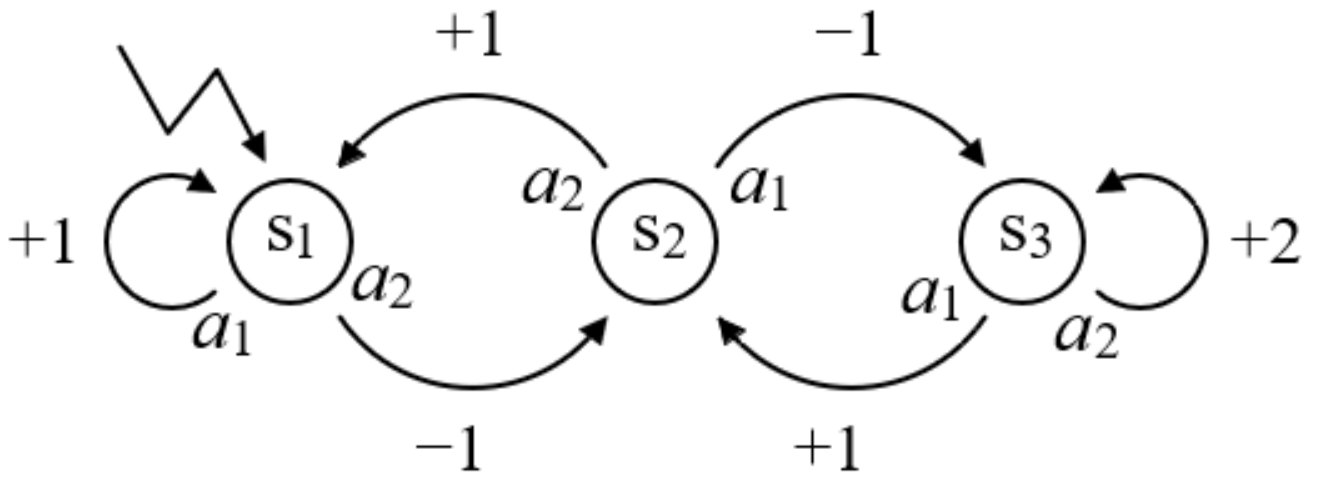
You may have noticed that we make the same kind of judgements about exploration versus exploitation in everyday situations. Should you eat at the old restaurant you have visited many times, or should you try the newly opened restaurant across the street which might be worse but might be much better?

Optional video



Delayed Reinforcement

The need for exploration also applies in the general case where there are multiple states, and where we may need to take a whole sequence of actions in order to get to the state from which we can obtain a high reward.



Recall that every policy π determines a Value Function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ where $V^\pi(s)$ is the expected discounted reward received by an agent who begins in state s and chooses its actions according to policy π .

If $\pi = \pi^*$ is optimal, then $V^*(s) = V^{\pi^*}(s)$ is the maximum (expected) discounted reward obtainable from state s . Knowing this optimal value function can help to determine the optimal policy.

Computing the Value Function

Let $\gamma = 0.9$ and consider the policy $\pi : S_1 \mapsto a_2, S_2 \mapsto a_1, S_3 \mapsto a_2$. We have

$$\begin{aligned}
 V^\pi(S_3) &= \frac{2}{1 - \gamma} = 20 \\
 V^\pi(S_2) &= -1 + \gamma V(S_3) = -1 + 18 = 17 \\
 V^\pi(S_1) &= -1 + \gamma V(S_2) = -1 + 15.3 = 14.3
 \end{aligned}$$

For this example, the policy π shown above must be the **optimal** policy, because the value function for states S_1 and S_2 under any other policy would not be more than 10. Therefore, the optimal value function is

$$V^* : S_1 \mapsto 14.3, S_2 \mapsto 17, S_3 \mapsto 20$$

If we were given this value function V^* , we could use it to determine the optimal policy π^* provided we also know the reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and the transfer function $\delta : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$. This information is sometimes called the "World Model".

Q-Function

The Q-function is a more sophisticated version of the value function which enables an agent to act optimally without needing to know the World Model.

For any policy π , the **Q-function** $Q^\pi(s, a)$ is the expected discounted reward received by an agent who begins in state s , first performs action a and then follows policy π for all subsequent timesteps.

If $\pi = \pi^*$ is optimal, then $Q^*(s, a) = Q^{\pi^*}(s, a)$ is the maximum (expected) discounted reward obtainable from s , if the agent is forced to take action a in the first timestep but can act optimally thereafter.

If the optimal Q-function Q^* is known, then the optimal policy is given by

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

The Q-function for the environment shown above can be computed as follows:

$$Q(S_1, a_1) = 1 + \gamma V(S_1) = 1 + 0.9 \times 14.3 = 13.87$$

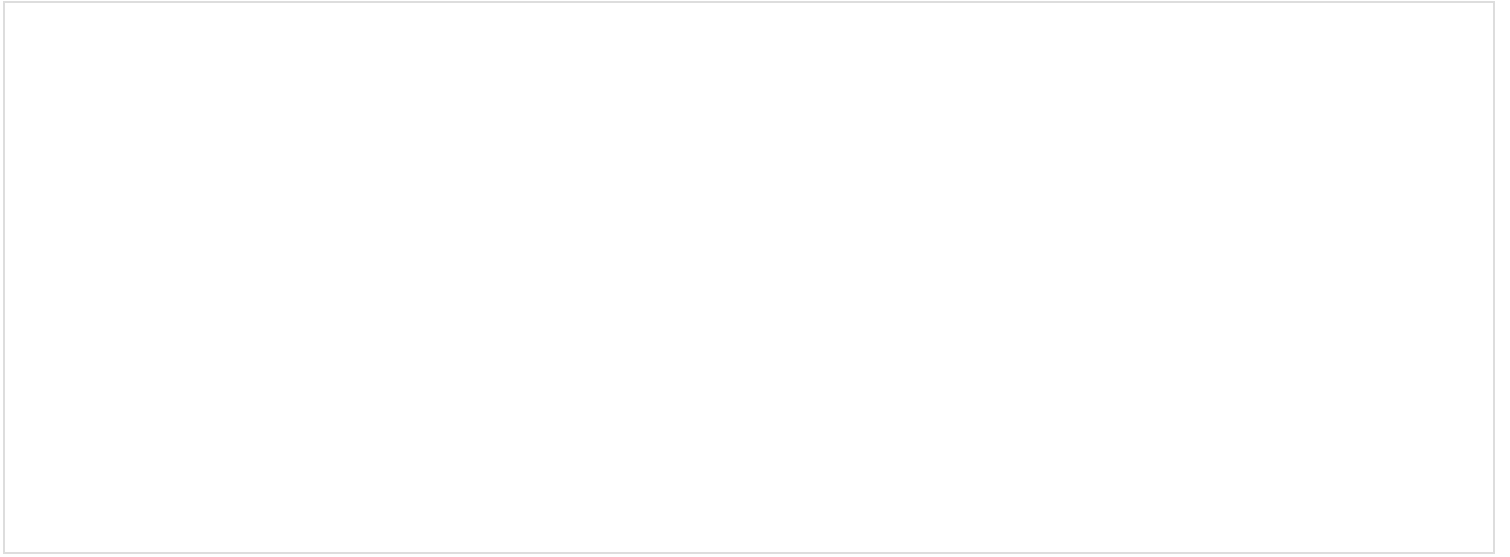
$$Q(S_1, a_2) = V(S_1) = 14.3$$

$$Q(S_2, a_1) = V(S_2) = 17$$

$$Q(S_2, a_2) = 1 + \gamma V(S_1) = 13.87$$

$$Q(S_3, a_1) = 1 + \gamma V(S_2) = 1 + 0.9 \times 17 = 16.3$$

$$Q(S_3, a_2) = V(S_3) = 20$$



TD-Learning and Q-Learning

Reinforcement Learning algorithms can generally be grouped into three classes:

1. **Value function learning**, including TD-Learning and Q-Learning
2. **Policy learning**, including policy gradients and evolution strategies
3. **Actor-Critic**, which is a combination of value function and policy learning

Value Function Learning

Recall that if $\pi = \pi^*$ is optimal, then $V^*(s) = V^{\pi^*}(s)$ is the maximum (expected) discounted reward obtainable from state s , and $Q^*(s, a) = Q^{\pi^*}(s, a)$ is the maximum (expected) discounted reward obtainable from s , if the agent is forced to take action a in the first timestep but can act optimally thereafter.

The idea behind value function learning is that the agent retains its own estimate $V()$ or $Q()$ of the "true" value function $V^*()$ or $Q^*()$. This estimate might be quite inaccurate to start with, but it gets iteratively improved over time so that it more closely approximates the true value. This process is sometimes called "Bootstrapping".

Temporal Difference Learning

Let's first assume that R and δ are deterministic. Then the (true) value $V^*(s)$ of the current state s should be equal to the immediate reward plus the discounted value of the next state

$$V^*(s) = R(s, a) + \gamma V^*(\delta(s, a))$$

We can turn this into an update rule for the estimated value:

$$V(s_t) \leftarrow r_t + \gamma V(s_{t+1})$$

If R and δ are stochastic (multi-valued), it is not safe to simply replace $V(s)$ with the expression on the right hand side. Instead, we move its value fractionally in this direction, proportional to a learning rate η .

$$V(s_t) \leftarrow V(s_t) + \eta [r_t + \gamma V(s_{t+1}) - V(s_t)]$$

Q-Learning

Q-learning is similar to TD-learning except that the Q-function $Q^* : S \times A \rightarrow \mathbf{R}$ depends on a state, action pair instead of just a state.

For a deterministic environment, π^* , Q^* and V^* are related by

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

$$Q^*(s, a) = R(s, a) + \gamma V^*(\delta(s, a))$$

$$V^*(s) = \max_b Q^*(s, b)$$

$$\text{So } Q^*(s, a) = R(s, a) + \gamma \max_b Q^*(s, b)$$

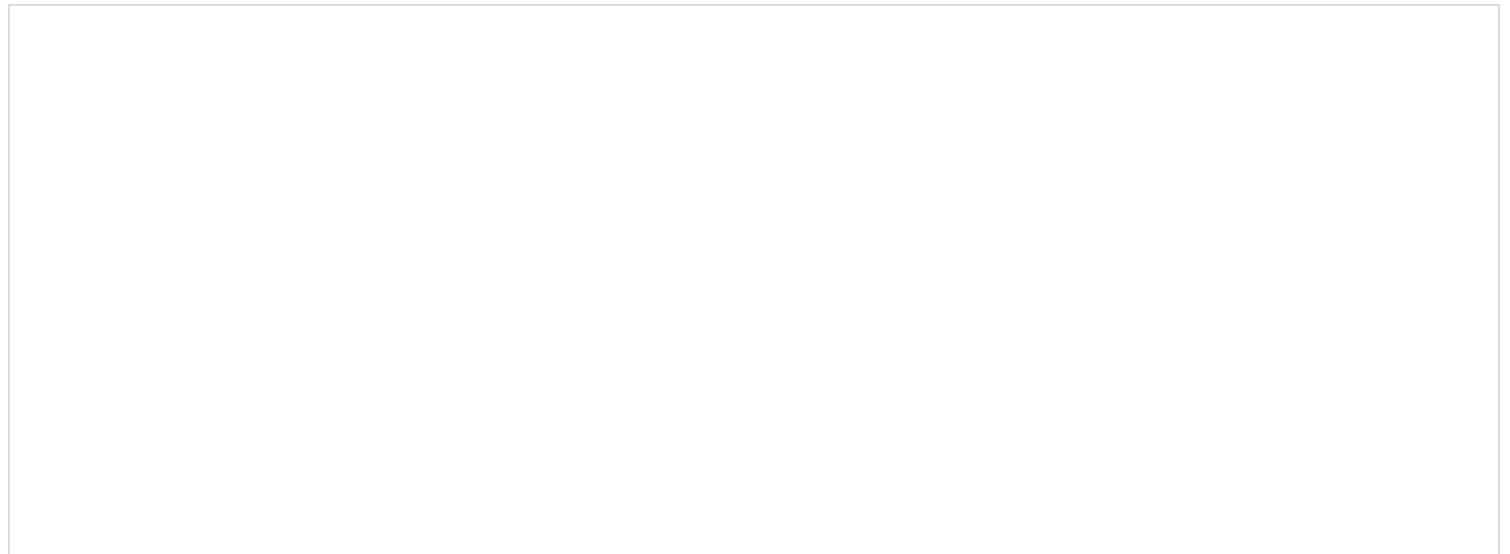
The last equation suggests that we can iteratively update our estimate Q by

$$Q(s_t, a_t) \leftarrow r_t + \gamma \max_b Q(s_{t+1}, b)$$

If the environment is stochastic, we instead write

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \eta [r_t + \gamma \max_b Q(s_{t+1}, b) - Q(s_t, a_t)]$$

Optional video



Theoretical Results and Generalization

It can be proved that TD-learning and Q-learning will eventually converge to the optimal policy, for any deterministic Markov decision process, assuming an appropriately randomised strategy ([Sutton, 1988](#); [Watkins & Dayan 1992](#); [Dayan & Sejnowski, 1994](#)).

However, these results rely on the assumption that every state will be visited many times. For small environments with a limited number of states and actions, this assumption is reasonable. For more challenging environments, the number of states can be exponentially large and it may be impractical

to assume that every state will be visited even once. In these situations, we need to rely on a model such as a neural network that is able to generalise to estimate the value of new states based on their similarity with states that have previously been encountered.

Optional video

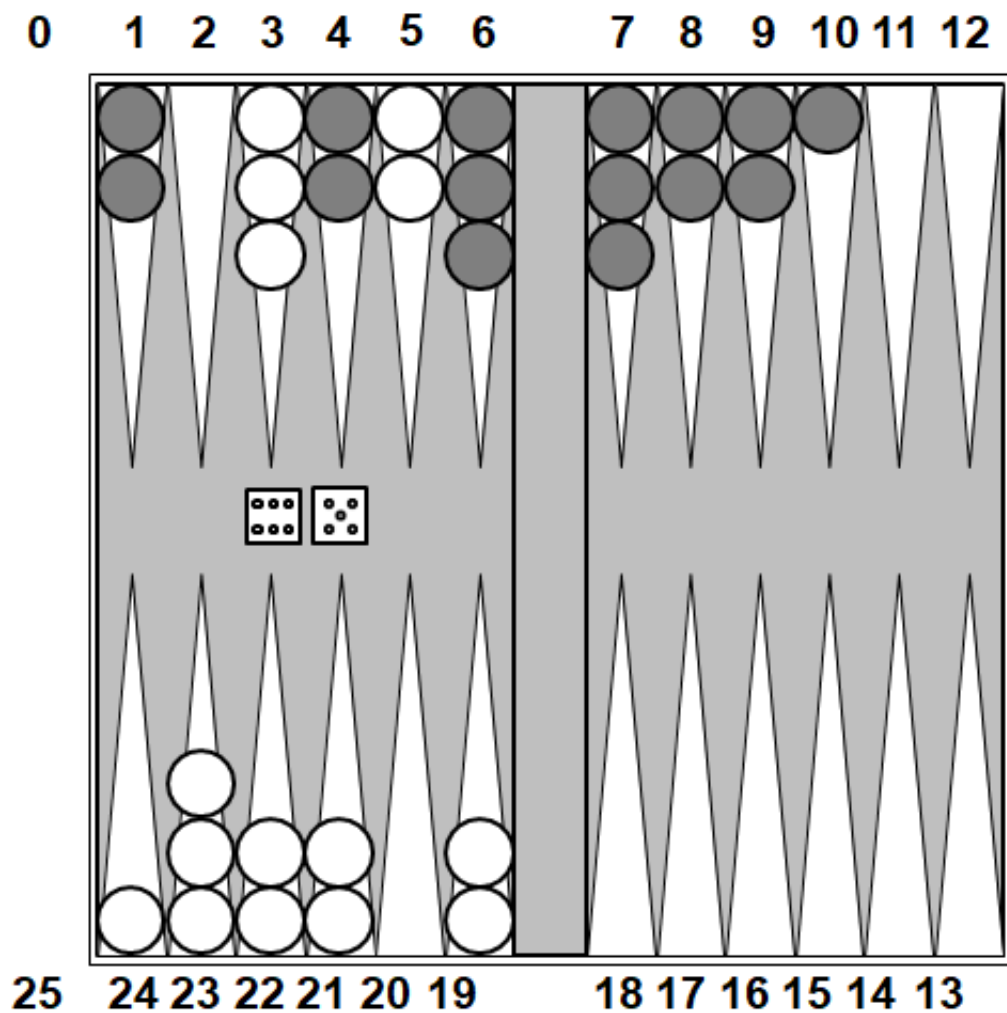


Computer Game Playing

Suppose we want to write a computer program to play a game like backgammon, chess, checkers or Go. This can be done using a tree search algorithm (expectimax, MCTS, or minimax with alpha-beta pruning). But we need:

- (a) an scheme for encoding any board position as a set of numbers, and
- (b) a way to train a neural network or other learning system to compute a board evaluation, based on those numbers.

Backgammon



One of the first practical applications of TD-learning (Tesauro, 1992) was to play the game of backgammon, which involves rolling dice and moving pieces around a board with 24 "points". In the above diagram, the white player is moving their pieces clockwise toward the bottom left quadrant while the black player is moving their pieces anti-clockwise toward the upper left quadrant. When a player gets all of their pieces into the last quadrant, they can use the dice rolls to move pieces off the board. The first player to get all of their pieces off the board is the winner.

If a player has only one piece on a point (such as point 10 or 24 above) and the other player moves a piece to that point, the original piece is sent back to the "bar" and must start again from the beginning. If a player has two or more pieces on a point, the other player is not allowed to move a piece to that point. We see in the above position that the black player has multiple pieces on points 1, 4, 6, 7, 8, 9 so that the white player will not be able to move their pieces on points 3 or 5 unless they roll either 2, 5 or 6.

TD-Gammon Neural Network

TD-Gammon used a 2-layer neural network with 196 inputs, 20 hidden nodes and 1 output. The inputs are arranged as follows:

- 4 units × 2 players × 24 points

- 2 units for pieces on the bar
- 2 units for pieces off the board

Four inputs are assigned to each colour and each point, with a 1-hot encoding used to specify one, two or three pieces and the fourth input scaled in proportion to the number of pieces if there are more than three. Additional inputs specify the number of pieces of each colour which are on the bar, or off the board.

For any encoded board state, the network produces an output V between 0 and 1 which is interpreted as its estimate of the probability of winning from that state.

At each timestep, the dice are rolled, the network considers all possible legal moves and computes the “next board position” that would result from each move. These are converted to the appropriate input format, fed to the network, and the one which produces the largest output is chosen.

The network can be trained by backpropagation using this formula

$$w_i \leftarrow w_i + \eta(T - V) \frac{\partial V}{\partial w_i}$$

where w_i are the weights in the network, η is the learning rate, V is the actual output of the network, and T is the target value.

The question is: How do we determine the target value T ? In other words, how do we know what the value of the current position “should” have been? Alternatively, how do we find a *better* estimate for the value of the current position?

How to choose the Target value

One approach is to have a human expert play many games and build a database of positions, dice rolls and chosen moves, similar to what was done for the ALVINN autonomous driving system we discussed in Week 1. The network could then be trained to increase the evaluation of each move chosen by the human expert and decrease that of other moves which were available but did not get chosen.

Another approach is for the network to learn from self-play by TD-learning, effectively using the evaluation of subsequent positions in the game to refine the evaluation of earlier positions.

Other methods such as TD-Root, TD-Leaf, MCTS and TreeStrap (Veness et al., 2009) combine learning with tree search. These are important for deterministic games like chess or Go, but less important for a stochastic game like backgammon because the randomness of the dice rolls limits the benefit of deep lookahead.

For backgammon, the agent receives just a single reward at the end of each game, which we can consider as the final value V_{m+1} (typically, +1 for a win or -1 for a loss). We then have a sequence of

game positions, each with its own (estimated) value:

(current estimate) $V_t \rightarrow V_{t+1} \rightarrow \dots \rightarrow V_m \rightarrow V_{m+1}$ (final result)

In this context, TD-learning simplifies and becomes equivalent to using the value of the next state (V_{t+1}) as the training value for the current state (V_k). A fancier version, called TD(λ), uses a weighted average T_k over future estimates as the training value for V_k , where

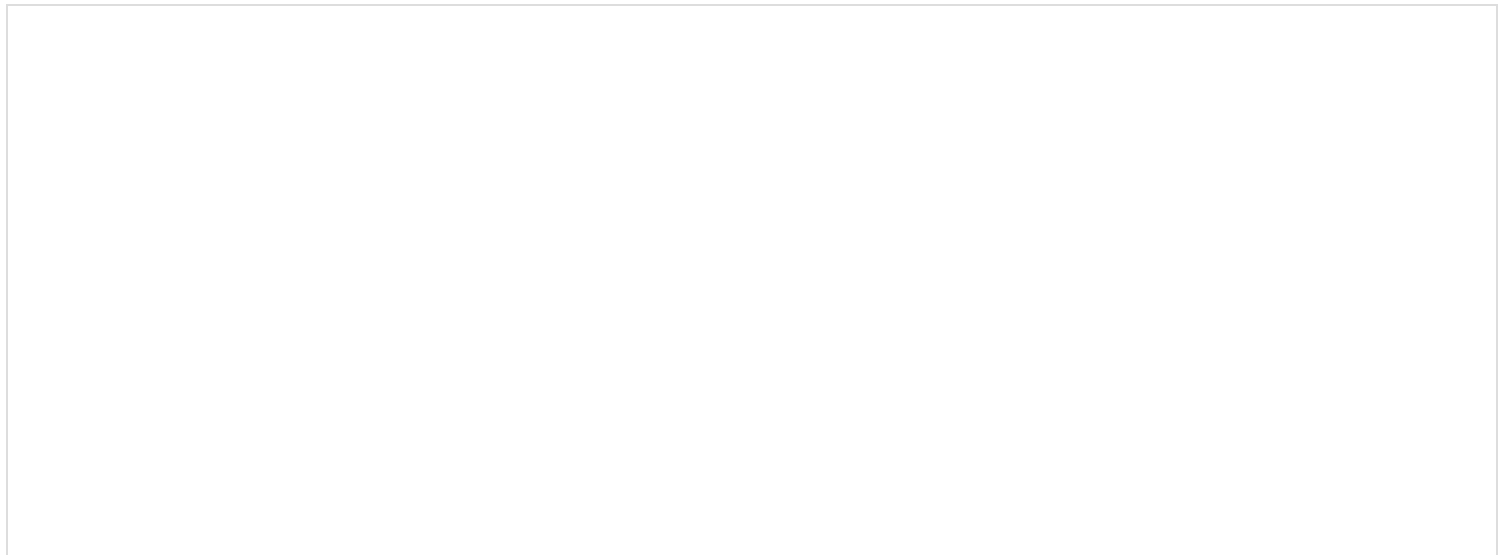
$$T_k = (1 - \lambda) \sum_{k=t+1}^m \lambda^{k-1-t} V_k + \lambda^{m-t} V_{m+1}$$

The parameter λ between 0 and 1 serves as a kind of discount factor, but is different to the usual discount factor γ (which can be equal to 1 in the case of TD-Gammon, because we do not care when we get the reward so long as we get it).

TD-Gammon

Tesauro trained two networks — one using Expert Preferences (EP), the other by TD-learning (TD). The TD-network outperformed the EP-network, and with modifications such as 3-step lookahead (expectimax) and additional hand-crafted input features, became the best backgammon player in the world in 1995.

Optional video



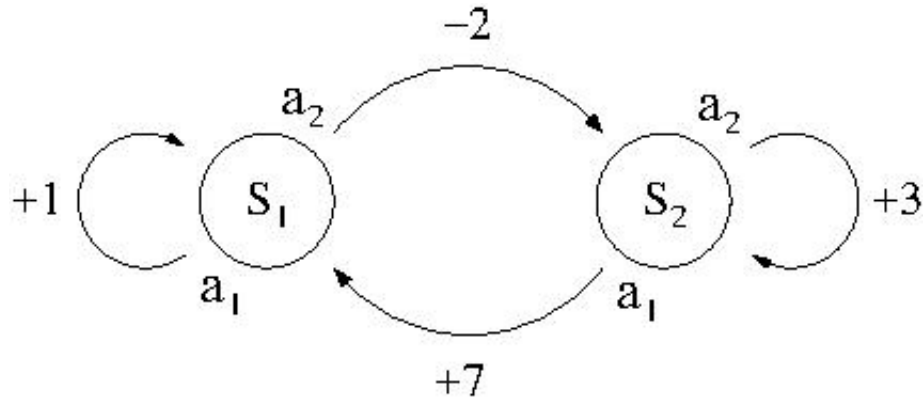
References

- Sutton, R.S. 1988. [Learning to predict by the methods of temporal differences](#), *Machine Learning*, 3(1), 9-44.
- Watkins, C.J. 1989. [Learning from Delayed Rewards](#), PhD Thesis, University of Cambridge.
- Watkins, C.J., & Dayan, P. 1992. [Q-learning](#), *Machine Learning*, 8(3-4), 279-292.

- Dayan, P., & Sejnowski, T.J. 1994. [TD \(\$\lambda\$ \) converges with probability 1](#), *Machine Learning*, 14(3), 295-301.
- Tesauro, G., 1992. [Practical issues in temporal difference learning](#). *Machine Learning*, 8(3), 257-277.
- Veness, J., Silver, D., Uther, W. & Blair, A., 2009. [Bootstrapping from game tree search](#), *Advances in Neural Information Processing Systems* (NIPS 22), 1937-1945.
-

Exercise: Reinforcement Learning

Consider an environment with two states $S = \{S_1, S_2\}$ and two actions $A = \{a_1, a_2\}$ where the (deterministic) transitions δ and reward R for each state and action are as follows:



$$\delta(S_1, a_1) = s_1, \quad R(S_1, a_1) = +1$$

$$\delta(S_1, a_2) = s_2, \quad R(S_1, a_2) = -2$$

$$\delta(S_2, a_1) = s_1, \quad R(S_2, a_1) = +7$$

$$\delta(S_2, a_2) = s_2, \quad R(S_2, a_2) = +3$$

Question 1 Submitted Oct 10th 2022 at 9:05:04 am

Assuming a discount factor of $\gamma = 0.7$, determine the optimal policy $\pi^* : S \rightarrow A$

asdf

Question 2 Submitted Oct 10th 2022 at 9:08:11 am

Still assuming $\gamma = 0.7$, determine the value function $V : S \rightarrow \mathbb{R}$

asdf

Question 3 Submitted Oct 10th 2022 at 9:14:17 am

Still assuming $\gamma = 0.7$, determine the values of the Q-function $Q : S \times A \rightarrow \mathbb{R}$

asdf

Question 4 Submitted Oct 10th 2022 at 9:17:25 am

Still assuming $\gamma = 0.7$, trace through the first few steps of the Q-learning algorithm, assuming a learning rate of 1 and with all Q values initially set to zero. Explain why it is necessary to force exploration through probabilistic choice of actions, in order to ensure convergence to the true Q values.

Here are some hints to get you started:

Since the learning rate is 1 (and the environment deterministic) we can use this Q-Learning update rule:

$$Q(S, a) \leftarrow r(S, a) + \gamma \max_b Q(\delta(S, a), b)$$

Let's assume the agent starts in state S_1 . Because the initial Q values are all zero, the first action must be chosen randomly. If action a_1 is chosen, the agent will get a reward of +1 and the update will be

$$Q(S_1, a_1) \leftarrow 1 + \gamma \times 0 = 1$$

asdf

Question 5 Submitted Oct 10th 2022 at 10:01:39 am

Now let's consider how the value function changes as the discount factor γ varies between 0 and 1. There are four deterministic policies for this environment, which can be written as π_{11} , π_{12} , π_{21} and π_{22} , where $\pi_{ij}(S_1) = a_i$, $\pi_{ij}(S_2) = a_j$

Calculate the value function $V_{(\gamma)}^{\pi} : S \rightarrow R$ for each of these four policies (keeping γ as a variable)

 5a_5.png

Question 6 Submitted Oct 10th 2022 at 10:01:01 am

Determine for which range of values of γ each of the policies π_{11} , π_{12} , π_{21} , π_{22} is optimal.

 5a_6.png

 5a_6b.png

Revision 7: Reinforcement Learning

This is a revision quiz to test your understanding of the material from Week 5 on Reinforcement Learning.

You must attempt to answer each question yourself, before looking at the sample answer.

Question 1 *Submitted Oct 10th 2022 at 10:05:22 am*

Describe the elements (sets and functions) that are needed to give a formal description of a reinforcement learning environment. What is the difference between a deterministic environment and a stochastic environment?

asdf

Question 2 *Submitted Oct 10th 2022 at 10:21:19 am*

Name three different models of optimality in reinforcement learning, and give a formula for calculating each one.



5a_2.png

Question 3 *Submitted Oct 10th 2022 at 10:33:51 am*

What is the definition of:

- the optimal policy
- the value function
- the Q-function?



5a_3.png

Question 4 *Submitted Oct 10th 2022 at 10:54:54 am*

Assuming a stochastic environment, discount factor γ and learning rate of η , write the equation for

- Temporal Difference learning TD(0)
- Q-Learning

Remember to define any symbols you use.
