

Super Simple Q-Learning Game

from "Teaching an AI to play a simple game using Q Learning" by Joren D

The game is the player 'P' must move to catch the cheese 'C' without falling into the pit 'O'

O = = P = < = = = = C = #
0 1 2 3 4 5 6 7 8 9 10 11 12

- Our strategy will be to make a program that allows the user to play the game, next we will add a random play mode, and finally we will add in the Q learning algorithm
- As previously discussed Q learning is a reinforcement learning algorithm. It works by allowing an agent to learn to reach a goal by taking random actions, however actions that bring the agent closer to goal are rewarded.

① From an implementation point of view, Q learning works by keeping a table of all game states, and the actions available at each of these states. For each state S and action A pair, a numerical value Q representing the possible reward for taking action A at state S is kept

- The Q table is initialized with random values.
- The Q table is updated over time using an algorithm that balances exploration of the state space with use of the table (to realize the agent's goals)

The table is updated using the equation below

$$Q(S_{t+1}, a_{t+1}) = Q(S_t, a_t) + \alpha \cdot (r_t + \gamma \cdot \max_a Q(S_{t+1}, a) - Q(S_t, a_t))$$

Diagram labels for the equation above:

- $Q(S_{t+1}, a_{t+1})$: new value
- $Q(S_t, a_t)$: old value
- α : learning rate
- r_t : instantaneous reward
- γ : discount factor
- $\max_a Q(S_{t+1}, a)$: new reward available by taking best action from state S_{t+1} which lands agent in S_{t+1}
- $- Q(S_t, a_t)$: old value

at S_t we take optimal action, we land in S_{t+1} giving agent max reward possible at that state

new reward available by taking best action from state S_{t+1} which lands agent in S_{t+1}

(2)

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha (r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

- this looks different from what we have seen in Mitchell's book, but similar to what is in Sutton & Barto
- to clarify, let $\alpha = 1$

$$Q(s_t, a_t) = Q(s_t, a_t) + 1 \cdot (r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$$

$$= \cancel{Q(s_t, a_t)} + r_t + \gamma \max_a Q(s_{t+1}, a) - \cancel{Q(s_t, a_t)}$$

$$Q(s_t, a_t) = r_t + \gamma \max_a Q(s_{t+1}, a)$$

that is what is seen in Mitchell

$$Q(s_t, a_t) = \overset{\text{so}}{r_t} + \gamma * \max_a (Q(s_{t+1}, a))$$

new value = $\underset{\substack{\uparrow \\ \text{immediate} \\ \text{reward}}}{r_t} + \text{discounted best available action}$

③ The steps of the algorithm are:

- Initialize the Q table with random values / or 0
for each s_t in S
for each a_t available at s_t

$$Q(s_t, a_t) = \text{random} / \text{or } 0$$

- While (playing) {
 $\text{chance} = \text{random}(1, 100)$
 if ($\text{chance} < \text{threshold}$) {
 select random action - a_t
 }
 else {
 select a_t such that it is
 max a_t for $Q(s_t, a_0 \dots a_m)$
 }

do action a_t , meaning

$$\text{next State} = \underset{\uparrow}{Q} \text{ Matrix}(a_t)$$

similar to adjacency matrix,
here we have a table that tells
us what the next state given
an action a_t at the current
state is.

Update Q table

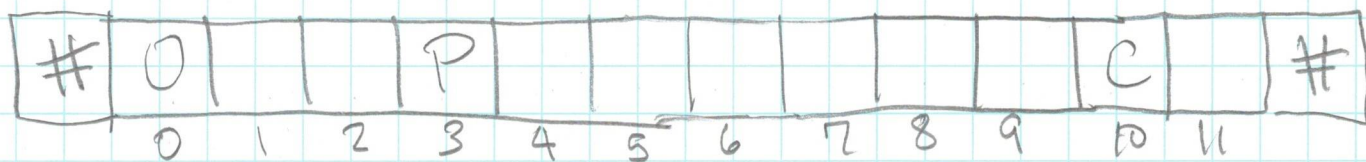
$$Q(s_t, a_t) = r_t + \max_{\text{State}} Q(\text{next}, a_i)$$

$$\} \quad s_t = \text{next state}$$

a_i is highest
reward a at
next state

④

The Q-matrix or Qadj identifies the action / actions needed to move from state to state



going left
will not take P
off game board

left right here is the game board

Q Matrix = $\begin{bmatrix} [0, 1], [0, 2], [1, 3], [2, 4], [3, 5], \\ [4, 6], [5, 7], [6, 8], [7, 9], [8, 10], \\ [9, 11], [10, 11] \end{bmatrix}$

state

in last state, going right
will not take you off the
game board

* As an aside, in graph learning, the game states are not known. It is up to the exploration to find them.

- In q learning the states are known, the actions that move the agent from one state to another are known, they are recorded in Q-matrix, or Q-adj.
- The goal of q learning is to find the most efficient way to get to the goal state.

⑦ Generally in Q-learning / reinforcement learning the Q table is used to train a neural network.

