

# Reinforcement Learning

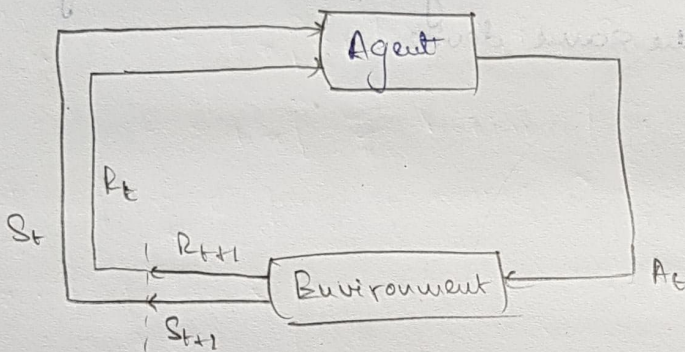
## (Deep Q-Learning)

Markov decision process (MDP): It is a way to formalize sequential decision making.

Components of MDP:

- agent
- environment
- state
- action
- rewards

\* Agent wants to maximize the cumulative reward.



Steps:

- At  $t$ , env. is in state  $S_t$ .
- Agent observes the state and takes action  $A_t$ .
- Env changes to  $S_{t+1}$  with a reward  $R_{t+1}$  is given to agent.
- Now,  $t+1$  is the current time step.

$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \dots$

$S_t, S_{t+1} \in S$   
 $A_t, A_{t+1} \in A$   
 $R_t, R_{t+1} \in R$

finite sets

finite MDPs

If infinite, then Infinite MDPs.

Hence,  $S_t, R_t$  (random vars) will have well defined prob. dist. (every value of  $S_t, R_t$  have associated prob.).

The prob. dist of  $S_t$  and  $R_t$  depend on the prev. time step (action and state)

\* Probability of transition to state  $s'$  with reward  $r$  from taking action  $a$  in state  $s$ .

$$p(s', r | s, a) = Pr \{ S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a \}$$

Expected Return:

\* Sum of future rewards.

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

$T \rightarrow$  final time step

If  $T \rightarrow$  finite, then the agent environment interaction can be broken into sub-sequences called episodes. Each episode ends in a terminal state  $T$ .



and the env. is reset after that. Next episode is independent of how prev. one ended. These are called as episodic tasks.

$G_t$  will be finite if ~~at~~ each reward is finite.

If  $T = \infty$

$$G_t = R_{t+1} + R_{t+2} + \dots, \text{ here } G_t = \infty.$$

Continuing tasks (can't be broken into episodes)

### Discounted Returns:

Future rewards are more heavily discounted, so the agent cares more about the immediate reward.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \quad 0 \leq \gamma \leq 1$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

If  $\gamma = 0 \Rightarrow G_t = R_{t+1}$  (the agent will choose greedy action to optimize the immediate reward)

\* Smaller the  $\gamma$   $\Rightarrow$  more myopic is the agent

Now: if  $R_t < \infty$  and  $0 < \gamma < 1 \Rightarrow G_t < \infty$

$$\text{Also: } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= R_{t+1} + \gamma [R_{t+2} + \gamma R_{t+3} + \dots]$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

Policies (what is the prob. of the agent to select any given action in a particular state).

For each state  $s \in S$ ,  $\pi$  is a prob. dist over  $a \in A(s)$ .

$\pi(a|s)$ : Prob. of taking action  $a$  in state  $s$ .

Value Functions (how good is it for an agent to be in a state or how good is it for the agent to choose an action in a given state).

Value functions depends on actions taken and hence on policies and thus on expected return.



• State value func (how good is any state for an agent following policy  $\pi$ )

$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] \\ = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

$V_{\pi}$ : Expected return from starting from state  $s$ , and following policy  $\pi$  thereafter.

• Action value func (how good it is for an agent to take any given action in a given state when following policy  $\pi$ )

Q-func

$$Q_{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a] \\ = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

\* Output of Q-func. is called Q-value.

$Q_{\pi}$ : Expected return when starting from state  $s$ , taking action  $a$  and then following  $\pi$  thereafter.

Optimality:

RL algorithms seek to find a policy that will yield more return to the agent if it follows that policy.

Optimal Policy:

$$\pi \geq \pi' \text{ if and only if } V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in \mathcal{S}.$$

Policy  $\pi$  is better than  $\pi'$

\* A policy  $\pi$  is better than  $\pi'$  if and only if, the state value functions for all states under policy  $\pi$  is greater than that under  $\pi'$ .

Optimal State value func:

The optimal policy has an associated optimal state-value func

$$V_*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s \in \mathcal{S}.$$

maximise  $V_{\pi}(s)$  over all possible  $\pi$



### Optimal Action Value Func:

optimal policy also has optimal action value func ( $Q$  func.)

$$Q_*(s, a) = \max_x Q_x(s, a) \quad \forall s \in S, a \in A(s)$$

### Bellman optimality eq. for $Q_*$ :

$Q_*$  must satisfy

$$Q_*(s, a) = E[R_{t+1} + \gamma \max_{a'} Q_*(s', a')]$$

The optimal action value func. of a state action pair is equal to the expected reward of choosing  $a$  in state  $s$  plus the max. expected discounted return from the next state action pair.

### Q-learning:

An RL algo that finds the optimal policy by learning optimal  $Q$ -values for each state action pair.

### Q-learning with value iteration:

Iteratively update  $Q$ -values for each state action pair using Bellman eqn. until the  $Q$ -func. converges to optimal  $Q$  func ( $Q_*$ ). This is called value iteration.

The  $Q$ -values are stored in  $Q$ -table.

Q-table	
states	actions
dim = [no. of actions $\times$ no. of states]	

### Exploration vs Exploitation:

$\epsilon$ -greedy strategy:  $\epsilon \rightarrow$  exploration rate

Initially  $\epsilon = 1$  and is reduced after every time step towards 0.

At every time step, generate a random number  $x$  between 0 and 1.

If  $x < \epsilon \rightarrow$  explore

$x > \epsilon \rightarrow$  exploit (become greedy)

(choose action with highest  $Q$ -value)



## Updating the Q-value:

We want the Q-value to eventually converge to optimal Q-value.

optimal Q-value    old Q-value

$$Q_*(s,a) - Q(s,a) = \text{loss}$$

$$= E \left[ R_{t+1} + \gamma \max_{a'} Q_*(s', a') \right] - E \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right] = \text{loss}$$

We update the Q-value of a state-action pair iteratively to minimize the loss whenever we encounter the same state action pair.

## Learning Rate ( $\alpha$ ):

A number between 0 and 1 which tells how quickly agent abandons the prev Q-value in the Q-table for the new value.

$\alpha = 1 \Rightarrow$  prev value is totally discarded.

$$\underbrace{Q^{\text{new}}(s,a)}_{\text{new Q-value}} = (1-\alpha) \underbrace{Q(s,a)}_{\text{old value}} + \alpha \underbrace{(R_{t+1} + \gamma \max_{a'} Q(s', a'))}_{\text{learned value}}$$

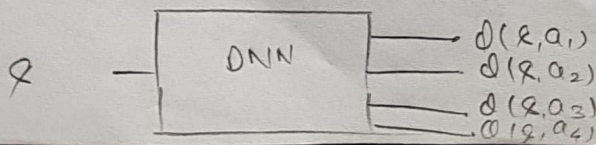
## Creating Lizard Game:

+1 ③	-1 ①	+5 ②
-1 ③	-10 ④	-1 ⑤
5 ⑥	-1 ⑦	+10 ⑧

## Deep Q-learning:

When state-spaces increase in size, iterative approaches such as value iteration become computationally expensive. Instead of using value iteration to compute optimal Q-function, we can use neural networks to approximate the optimal Q-function.

\* The DNN accepts state of the env. as input and outputs estimated Q-values for each action in that state. The DNN tries to approximate the optimal Q-function that satisfies Bellman's eqn.





\* Loss is calculated by comparing the output  $Q$ -values of the network to the target  $Q$ -values from the RHS of Bellman Eqn.

Based on the loss, weights of the network are updated using backprop. The process is repeated over and over for each state until loss is minimized.

\* Consecutive frames are often used as input to the NN after some preprocessing.

\* Activation func. are not added at the output layer since we want raw  $Q$ -values from the network.

### Replay Memory:

We store the agent's experiences at each time step into a dataset called replay memory.

$$x_t = (s_t, a_t, r_{t+1}, s_{t+1})$$

↳ agent's experience at time step  $t$ .

Last  $N$  experiences are stored in replay memory (finite size)

\* We randomly sample from replay memory to train the DQN and not train from sequential experience to break the correlation between consecutive samples. (for more efficient learning)

Deep  $Q$ -Network is also called the policy network.

### Training a Policy Network:

Assuming that we input a single state as input (batch size = 1)

For eg., we pass  $s_4$  as input.

$$x_4 = (s_4, a_4, r_5, s_5)$$

The network outputs  $Q$ -values for all the action we select the  $Q$ -value for the action  $a_4$  from the exp. tuple. Loss is then calculated by comparing this <sup>estimated</sup>  $Q$ -value with the target  $Q$ -value for  $a_4$ .

Now, target  $Q$ -value is given by the RHS of Bellman eqn.

$$Q_*(s_4, a_4) = R_{t+1} + \gamma \max_{a_5} Q_*(s_5, a_5)$$

To compute  $\max_{a_5} Q_*(s_5, a_5)$ , we pass  $s_5$  to the policy network and choose the max  $Q$ -value output over all the actions.



Therefore, two forward passes are done before calculating the loss and updating the weights.

\*In case of value iteration,  $\max_{a'} q_*(s', a')$  was found directly from the Q-table.

### Training Issues with DQNs:

The first pass through DQN (Policy Network) is done to find the estimated Q-value and 2nd pass for target Q-value. Both of these passes are done through the same network before updating the weights.  $q(s, a)$  approaches  $q_*(s, a)$  before updating after updating the weights but  $q_*(s, a)$  also moves away. This introduces instability during training.

### The Target Network:

A clone of policy net called target network is used to do the 2nd pass to find the target Q-value. The weights of this network is updated to the same as policy network after some time steps. This reduces instability.