

Chapter 11

Reinforcement Learning

Contents

11 Reinforcement Learning	1
11.1 TD Learning	2
11.2 SARSA	2
11.3 Q-Learning	3
11.4 Deep Q-Learning	3
11.5 Policy Gradient Methods	4
11.6 Actor-Critic	4

11.1 TD Learning

Algorithm 11.1 Tabular TD(0) for estimating v_π

```
1: function TD( $\pi$ : the policy to be evaluated)
2:   Initialize  $V(s)$  arbitrarily (e.g.,  $V(s) = 0, \forall s \in S^+$ )
3: end function
```

- Model-free
- Off-policy
- Discrete action and state spaces

11.2 SARSA

- SARSA works by learning a state-action value function rather than a state value function like TD learning.
- SARSA explore the transitions from one state-action pair to another state-action pair and learns the state-action value function.
- In on-policy learning, the optimal value function is learned from actions taken using the current policy.

- **Update** rule:

$$Q(S_t, A_t) = \quad (11.1)$$

- As in all on-policy methods, Q_n is continually esimated for the policy π , while the policy π is updated using an ϵ -greedy policy.

Algorithm 11.2 ϵ -Greedy Policy

```

1:  $p \leftarrow \text{RANDOM}$ 
2: if  $p < \epsilon$  then
3:   pull random action
4: else
5:   pull current-best action
6: end if
```

Algorithm 11.3 SARSA (on-policy TD control)

```

1: Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\epsilon > 0$ 
2: Initialize  $Q(s, a)$ , for all  $s \in S^+$ ,  $a \in A(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$ 
3: loopfor each episode:
4:   Initialize  $S$ 
5:   Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
6: end loop
```

- Model-free.
- On-policy.
- Discrete action and state spaces.

11.3 Q -Learning

- Q -Learning is very similar to SARSA. The major difference lies in the update rule of the Q -function.

11.4 Deep Q -Learning

- For tasks with continuous state space, updating the Q -function for all state-action pairs can be computationally inefficient and infeasible.
- Rather than using value learning to directly find the optimal Q -function, a function estimator can be used to estimate the optimal Q -function.

- ANNs are effective function estimators. Deep neural network (DNN) can be used to estimate the Q -function for each state-action pairs.
- DQN is trained using batch stochastic gradient updates and experience replay. Experience replay can interact with the environment to generate training data for the DQN.
- Experience replay is a technique where the agent stores a subset of its experiences $\langle s, a, r, s' \rangle$ in a memory buffer and samples from this buffer to update the Q -function.
- Experience replay selects an ϵ -greedy action from the current state, executes it in the environment, and gets back a reward and the next state.
- If DQN was trained with single samples, each sample and the corresponding gradients
- ...
- Model-free.
- Off-policy.
- Continuous state space.
- Discrete action space.

11.5 Policy Gradient Methods

- Policy gradient methods learn a parameterized policy than can select actions without consulting a value function. The parameters of the policy are called policy weights.
- Policy gradient methods are methods for learning the policy weights using the gradient of some performance measure with respect to the policy weights.
- Policy gradient methods seek to maximize performance and so the policy weights are updated using gradient ascent.

11.6 Actor-Critic

- Actor-Critic method is a TD version of policy gradient. It uses two neural networks, one actor network and one critic network.
- The actor network decides what action should be taken and the critic network informs the actor network how good was the action and how it should update to improve.
- The learning of the actor network is based on policy gradient approach.
-