

# Report of MVIG Training, Deep Reinforcement Learning

Yongxi Huang  
SJTU

huangyongxi@sjtu.edu.cn

## Abstract

*Reinforcement learning is an important topic of machine learning. By reading **Reinforcement Learning: An Introduction** and learning from the online lecture by David Silver, I have learned some concepts and implementation of reinforcement learning. In this report, I want to do a summary and combining of the knowledge learned. I also have trained some models. However, because of the limited time, I can't trained them totally. But I want give some result of the training in this report. And I'm still working on the code these days.*

## 1. Introduction

To do a summary of the knowledge, first I want to combining some basic concept of reinforcement learning. And then introduce some classical algorithms solving RL problem such as DP, MC and Q-learning. Finally, I want to give my understanding of some approaches introduced in some paper. However, I may have some misunderstanding about RL or make some mistakes. I hope anyone reading this report can correct them for me.

## 2. Basic concepts

### 2.1. Reinforcement Learning

Reinforcement learning is a kind of unsupervised learning method. In deep-learning, if we want to train the AI to do something, we must give it some correct answer like annotations and labels, and it can learning from the data by minimize the loss function. That is some kind of supervised learning, which means learning from some given labeled result. It have some feedback indicates the correct action to take, independently of the action actually taken

But in reinforcement learning, there is no supervisor. The agent can learn only from the reward signal it obtains. Besides, the reward feedback can be delayed and the actions of agents can change the environment and the reward. In summary, reinforcement learning agents learn by trail-and-error.

Reinforcement learning cant work on any kinds of environments. The agents always learn from its current state and the predicted future state. That means the current state must completely characterizes the process, otherwise the history of the process may make the learning process costly. In other words, the process of a reinforcement learning must has the Markov property:

$$\mathbb{P} = [S_{t+1}|S_t] = \mathbb{P} = [S_{t+1}|S_1, S_2, \dots, S_t]$$

In short, the Markov property means the future is independent of the past given the present.

The environment with Markov states that the reinforcement learning agents work on is called Markov Decision Process (MDP). Its a Markov process with reward and actions. In fact, in Reinforcement Learning: An Introduction , reinforcement learning is regarded as any method that is suited to solving MDP problem. A MDP can be decision as  $\langle S, A, P, R, \gamma \rangle$ . There are several useful concept about MDP.

### 2.2. States

In my opinion, this is a fundamental concept of MDP. A state contains all relevant information for the agent to make decision. However, when we talk about states, we are always talk about agent states instead of environment states. That means the states dont contains all the information of the environment, but the information that can be obtained by the agent. Besides, the states of a MDP have the Markov property. And the states can have different form.

### 2.3. Reward

The reward  $R_t$  is a scalar feedback based on the state or the action of the agent. The goal of the agent is to maximize the expected cumulative reward. The reward of a specific state can be defined as:

$$R_s = \mathbb{E}[R_{t+1}|S_t = s]$$

As we always care about the cumulative reward of a process, we can define the return of a process:

$$G_t = R_t + 1 + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

In that case, the goal of agent is to maximize the expected return.

## 2.4. Policy

Policy is a map from states of environment to the actions of agent. It define the behavior of an agent. It is defined as:

$$\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$$

Just as its definition, the policy  $\pi : S \rightarrow A$  has nothing to do with time. In practice, sometimes we estimate the future return according to the policy taken, which is called on-policy. And sometimes the policy doesn't influence the update of estimation of future return(policy evaluation), which is called off-policy.

## 2.5. Value function

Once we have reward and return, we can define the value function to evaluate a state or action. Its a prediction of future reward which can be defined as:

$$v_\pi = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

Value function is a function about states under the policy  $\pi$ . The agent choose different actions based on the value function. The value function defined above only evaluate the quality of states, which is called the state-value function. If we take the action further into account, we can define the action-value function:

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

It not only evaluate the quality of a state, but also the action taken on this state. In short, value function is the expected return. Find a proper way to evaluate value function is an vital part of reinforcement learning. After define these concept, the goal of reinforcement learning can be seen as find a optimal policy  $p_i^*$  that can maximize the value function.

$$\pi^* = \operatorname{argmax}_\pi V^\pi(s), \forall s$$

## 3. Basic Algorithms

### 3.1. Dynamic Programming

However, DP is not only a strategy to solve the reinforcement learning problem, but also an idea of other strategy. I think the reason why we can use Dynamic Programming to solve this problem is also the Markov property of MDP. Because the state-value function or action-value function have

nothing to do with the history of a process, we can decomposed them using the Bellman Expectation Equation:

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s] \\ q_\pi(s, a) &= \mathbb{E}_{\pi\pi}[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t = s, A_t = a] \end{aligned}$$

#### 3.1.1 Policy Evaluation

According to Bellman Expectation, the value function can be decomposed into immediate reward and the value of successor states. When given the policy  $\pi$ , we can calculate the value function of a specific state using Bellman Expectation Equation. We can see this process as a policy evaluation.

Using DP to evaluate a policy is like this:

---

#### Algorithm 1: Policy evaluation of DP

---

```

input : An policy  $\pi$ 
output: The estimated value function  $v_\pi$ 

1 while  $\Delta > \theta$  (a small positive number) do
2    $\Delta \leftarrow 0$ 
3   for each  $s \in S$  do
4      $temp \leftarrow v(s)$ 
5      $v(s) \leftarrow$ 
6        $\sum_a \pi(a|s) \sum_{s'} p(s'|s, a) [r(s, a, s') + \gamma v(s')]$ 
7      $\Delta \leftarrow \max(\Delta, |temp - v(s)|)$ 
7 Output  $v \approx v_\pi$ 

```

---

#### 3.1.2 Policy improvement

However, the algorithm above can only estimate the value function of a policy  $\pi$ . In other words, it can only do policy evaluation. If we want to find an optimal policy, we need to do policy improvement. As we can calculate the value function of an existing policy  $\pi$  according to the algorithm above, a simple way to find the optimal policy is to change the action of a specific state and see whether it can improve the value function. The action-value can be calculated by

$$\begin{aligned} q_\pi &= \mathbb{E}_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned}$$

If  $q_\pi(s, a)$  is greater than  $v_\pi(s)$ , its better to select  $a$  in  $s$ . and we can get a new greedy policy  $\pi'$  by

$$\begin{aligned} \pi'(s) &= \operatorname{argmax}_a q_\pi(s, a) \\ &= \operatorname{argmax}_a \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \operatorname{argmax}_a \sum_{s', r} p(s', r | a, a) [r + \gamma v_\pi(s')] \end{aligned}$$

This policy improvement algorithm is like:

---

**Algorithm 2: Policy iteration of DP**


---

**input :** Old policy  $\pi$ , value function  $V$ .  
**output:** Optimal policy

```

1 policy ← stable ← true
2 for each  $s \in S$  do
3   old ←  $\pi(s)$ 
4    $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r|s, a)[r + \gamma V(s')]$ 
5   if old ←  $\pi(s)$  then
6     policy ← stable ← false, break
7 if policy-stable then
8   return  $V \approx v_*$  and  $\pi \approx \pi_*$ 
9 else
10  go to policy evaluation

```

---

As we improve the policy iteratively in this algorithm, its called policy iteration. Another way to find the optimal policy is value iteration. We can see value iteration as combining the policy improvement and policy evaluation.

---

**Algorithm 3: Value iteration of DP**


---

```

1 Initialize array  $V$  arbitrarily
2 while  $\Delta > \theta$  (a small positive number) do
3    $\Delta \leftarrow 0$ 
4   for each  $s \in S$  do
5      $v \leftarrow V(s)$ 
6      $V(s) \leftarrow \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma V(s')]$ 
7      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
8 Output a deterministic policy,  $\pi \approx \pi_*$ , such that
9  $\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s', r|s, a)[r + \gamma V(s')]$ 

```

---

However DP is not really practical. First, the full backup operation used to update value requires high computation. And DP is based on all possible next states rather than on a sample next states, which means it is model based. So the classical DP algorithms are of limited utility in reinforcement learning both because of their assumption of a perfect model and because of their great computation expense. But DP is the base of other algorithms. An important concept created by DP, for example, is generalized policy iteration (GPI) which means letting policy evaluation and policy improvement processes interact. In **Reinforcement Learning: An Introduction**, it describe these two process as this:

One process takes the policy as given and performs some form of policy evaluation, changing the value function to be

more like the true value function for the policy. The other process takes the value function as given and performs some form of policy improvement, changing the policy to make it better, assuming that the value function is its value function.

And in my opinion, GPI is essential in many other methods solving reinforcement learning problems.

### 3.2. Monte Carlo Methods (MC)

A big improvement from DP to MC is that MC is model free which means we dont assume complete knowledge of the environment. However, model-free doesnt mean that there is no model. The model in MC just need to generate the sample transitions without the complete probability distributions. The basic idea of MC is value = mean return. And mean return can be estimated by empirical mean return. And the idea of GPI is just like DP. Besides, the value function can be updated incrementally.

$$N(S_t) \leftarrow N(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)}(G_t - V(S_t))$$

This formula uses a method that always used to estimate some target:

$$NewEstimate \leftarrow OldEstimate + StepSize[Target - OldEstimate]$$

The  $Target - OldEstimate$  is called as error, which is reduced by taking a step toward the Target. And when it comes to policy improvement, it proceeds according to GPI defined in DP. And no model is needed as the policy is made with respect to the current value function. But here is a dilemma between exploration and evaluation in MC method. To solve it, we can use  $\epsilon - greedy$  which means most of the time the agent chooses the action with maximal estimated action value, and with probability  $\epsilon$  it chooses other actions. A MC algorithm is shown in Algorithm 4.

The on-policy in fact makes a compromise that it learns action values for both optimal policy and the near-optimal policy. Another solution for the evaluation and exploration dilemma is off-policy learning. It use two different policy for evaluation (target policy) and exploration (behavior policy). As these two policy is separated, the target policy may be deterministic while the behavior policy can continue to sample all possible actions.

A drawback of MC is, as we can see, the updating of policy happens after each episode. So MC can only solve episodic tasks but not continuing tasks. An idea is to update the policy online. So here comes TD method.

## 4. TD learning

Temporal-Difference Learning (TD learning) can be seen as a combination of MC and DP. In MC, the value func-

---

**Algorithm 4: On-policy MC control**

---

```
1 Initialize
2  $Q(s, a) \leftarrow \text{arbitrary}$ 
3  $Returns(s, a) \leftarrow \text{emptylist}$ 
4  $\pi(a|s) \leftarrow \text{anarbitrary} \varepsilon - \text{softpolicy}$ 
5 while true do
6   Generate an episode using  $\pi$ 
7   for each pair  $s, a$  appearing in the episode do
8      $G \leftarrow$  following the first occurrence of  $s, a$ 
9     Append  $G$  to  $Returns(s, a)$ 
10     $Q(s, a) \leftarrow \text{average}(Returns(s, a))$ 
11   for each  $s$  in the episode do
12      $A_* \leftarrow \text{argmax}_a Q(s, a)$ 
13     for all  $a \in A(s)$  do
14        $\pi(a|s) = 1 - \varepsilon + \varepsilon/|A(s)|$ , if  $a = A_*$ 
15        $\pi(a|s) \varepsilon/|A(s)|$ , if  $a \neq A_*$ 
```

---

tion is updated until the return following the visit is known. The constant- MC for example, the increment is calculated through

$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

But TD methods need wait only until the next time step rather than the next episode. The value prediction of TD(0) is

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

The target for the TD update is  $R_{t+1} + \gamma V(S_{t+1})$  rather than  $G_t$ . This kind of updating method is called bootstrapping method. And its just what used in DP. In summary, MC uses a sample return as the real expected return is unknown. And DP use a target at the current estimate as the new estimate is unknown. The TD target combines DP and MC because it samples the expected values and it uses the current estimate  $V$  instead of the true  $v_\pi$ . A simple TD prediction algorithm TD(0) is shown as follow:

From another perspective, TD and MC both use sample backups while DP uses full backups. And we can see the update in TD(0) as updating by error  $[R + \gamma V(S') - V(S)]$  which measuring the difference between the estimated value of  $S_t$  and the better estimate  $R_{t+1} + V(S_{t+1})$ . This error is called TD error.

Another TD prediction method is batch updating which means although some states have been visited, the value function is changed only once by the sum of the increments after a batch of error. Many another RL algorithms use TD prediction.

---

**Algorithm 5: TD(0) for estimating  $v_\pi$** 

---

```
input : Policy  $\pi$ .
output: Value evaluation  $v_\pi$  of policy  $\pi$ 

1 Initialize  $V(s)$  arbitrarily
2 while true do
3   Generate an episode  $S$  using  $\pi$ 
4   for each step of episode do
5      $A \leftarrow \text{action from } S$ 
6     Take action  $A$  and observe  $R, S'$ 
7      $V(S) \leftarrow V(S) + \alpha[R + \gamma V(S') - V(S)]$ 
8      $S \leftarrow S'$ 
```

---

#### 4.1. Sarsa

Sarsa is an algorithm using TD prediction. In order improve the policy, it considers action value rather than state value, which means its on policy. Also, its a GPI process. In order to use TD(0) to update action value rather than state value, we can change the formula into this:  $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$ . A sarsa algorithm using  $\varepsilon - greedy$  is shown in Algorithm 6.

---

**Algorithm 6: Sarsa, On-policy TD control**

---

```
1 Initialize  $Q(s, a)$  arbitrarily
2 for each episode do
3   Initialize  $S$ 
4   Choose  $A$  from  $S$  using policy derived from  $Q$ 
   (e.g.  $\varepsilon - greedy$ )
5   for each step of episode do
6     Take action  $A$ , observe  $R, S'$ 
7     Choose  $A'$  from  $S'$  using policy derived from
        $Q$  (e.g.  $\varepsilon - greedy$ )
8      $Q(S, A) \leftarrow$ 
        $Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$ 
9      $S \leftarrow S'; A \leftarrow A'$ 
```

---

#### 4.2. Q-learning

Q-learning also use the idea of TD prediction, but its an off-policy TD control method. It is defined by

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

In that case, the updating of the action value function has nothing to do with the action taken before. So its off-policy. The algorithm is shown in Algorithm 7.

We can consider  $Q$  as the knowledge learned. And the optimal can be derived greedily from  $Q$ .

---

**Algorithm 7: Q-learning, Off-policy TD control**

---

```
1 Initialize Q(s,a) arbitrarily
2 for each episode do
3   Initialize S
4   for each step of episode do
5     Choose A from S using policy derived from Q
      (e.g.  $\epsilon - greedy$ )
6     Take action A, observe R, S'
7      $Q(S, A) \leftarrow$ 
       $Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ 
8      $S \leftarrow S'$ 
```

---

However, there is still many basic knowledge about reinforcement learning I don't understand as I just finish half of the reference book Reinforcement Learning: An Introduction. I will keep working on it. But for now, I want to talk about something I learn from some paper. However, it's really helpful for understanding the paper to learn some basic concept and fundamental algorithm of reinforcement learning.

## 5. Paper reading

### 5.1. DQN

In Playing Atari with Deep Reinforcement Learning, an approach using neural network is introduced. It uses the idea of Q-learning, while using convolutional neural network to estimate the reward. It is model free as we just input reward, terminal signals and possible actions. As the environment of the agent can't be described in a single screen shot, the method considers a sequence of actions and observations as one state. And the reward is depended on the scores of games.

Just as other reinforcement learning problem, we should do policy evaluation and policy improvement to find the optimal policy. For policy evaluation, a simple idea is using the value iteration just as classical Q-learning like  $Q_{i+1}(s; a) = E[r + \max_{a'} Q_i(s; a') | s, a]$ . But in this paper, it uses a function approximator to estimate the action value function. More specifically, it uses a neural network to approximate it, which is called Q-network. And the Q-network can be trained by minimizing the loss function  $L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$  using stochastic gradient descent, where  $y_i = \mathbb{E}_{s' \sim \epsilon} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ .

In that case,  $y_i$  can be regarded as the estimate value function and the loss function can measure the difference between the estimation and the knowledge it has learned. As for  $\rho(s; a)$ , I think it's a kind of behavior distribution like  $\epsilon - greedy$  method. However, the target  $y_i$  here is different from the target in deep learning, which is fixed before

training.  $y_i$  is changing in the training process.

### 5.2. Dueling Network

In Dueling Network Architectures for Deep Reinforcement Learning, it creates a new architecture of neural network that is suitable for RL.

## 6. Experiments

## 7. Results

## 8. Conclusion