

DQN 구현

[쉽게 구현하는 강화학습 2화]

판요랩 - 노승은, 전민영

2019.05.12

<https://github.com/seungeunrho/minimalRL>

복습 – (6장)

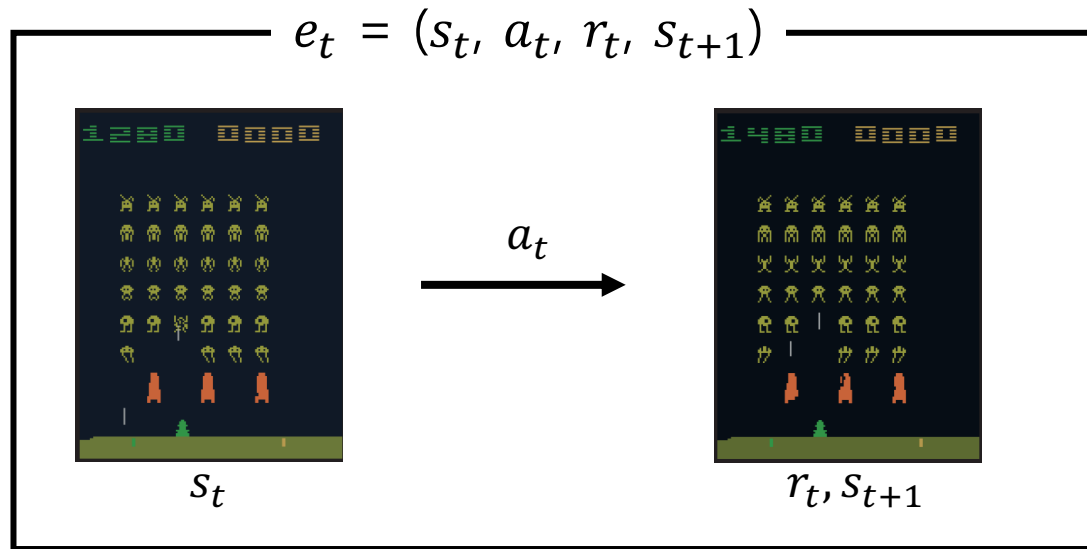
Value Function Approx. By Stochastic Gradient Descent

- Goal: find parameter vector \mathbf{w} minimising mean-squared error between approximate value fn $\hat{v}(s, \mathbf{w})$ and true value fn $v_\pi(s)$

$$J(\mathbf{w}) = \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, \mathbf{w}))^2]$$

- Have assumed true value function $v_\pi(s)$ given by supervisor
- In practice, we substitute a *target* for $v_\pi(s)$
- For TD(0), the target is the TD target $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$

복습 – Replay Buffer



$$D = \{e_1, e_2, e_3, \dots, e_N\}$$

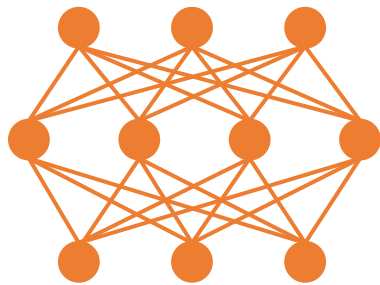


Replay Buffer

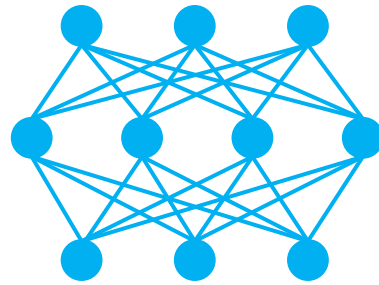
복습 – Target Network

$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(D)} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a; \theta_i) \right)^2 \right]$$

Target Network



Q Network



일정 주기마다 update
 $\theta_i^- = \theta_i$

Import

```
1 import gym
2 import collections
3 import random
4
5 import torch
6 import torch.nn as nn
7 import torch.nn.functional as F
8 import torch.optim as optim
```

- Collections library는 replay buffer에서 쓰일 deque 를 import 하기 위함.

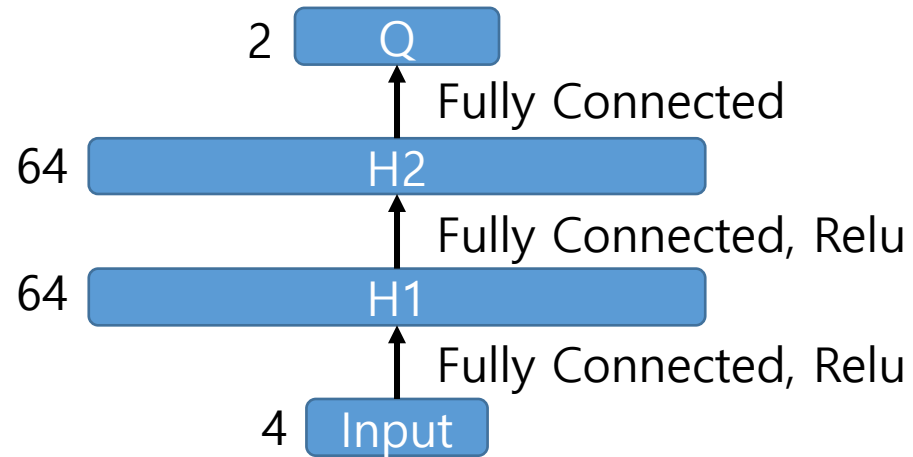
Replay Buffer

```
11 class ReplayBuffer():
12     def __init__(self):
13         self.buffer = collections.deque()
14         self.batch_size = 32
15         self.size_limit = 50000 ----->
16
17     def put(self, data):
18         self.buffer.append(data)
19         if len(self.buffer) > self.size_limit:
20             self.buffer.popleft()
21
22     def sample(self, n):
23         return random.sample(self.buffer, n)
24
25     def size(self):
26         return len(self.buffer)
27
```

Buffer의 최대 크기

Q Network

```
28 class Qnet(nn.Module):
29     def __init__(self):
30         super(Qnet, self).__init__()
31         self.fc1 = nn.Linear(4, 64)
32         self.fc2 = nn.Linear(64, 64)
33         self.fc3 = nn.Linear(64, 2)
34
35     def forward(self, x):
36         x = F.relu(self.fc1(x))
37         x = F.relu(self.fc2(x))
38         x = self.fc3(x)
39         return x
40
41     def sample_action(self, obs, epsilon):
42         out = self.forward(obs)
43         coin = random.random()
44         if coin < epsilon:
45             return random.randint(0,1)
46         else :
47             return out.argmax().item()
```



Main 1

```
72 def main():
73     env = gym.make('CartPole-v1')
74     q = Qnet()
75     q_target = Qnet()
76     q_target.load_state_dict(q.state_dict())
77     memory = ReplayBuffer()
78
79     avg_t = 0
80     gamma = 0.98
81     batch_size = 32
82     optimizer = optim.Adam(q.parameters(), lr=0.0005)
```

- Q와 Target Q 네트워크 두 개를 선언
- Q를 Target Q로 복사
- state_dict 는 model의 weight 정보를 dictionary형태로 담고 있음
- 예컨대 아래와 같음

Model's state_dict:

```
conv1.weight      torch.Size([6, 3, 5, 5])
conv1.bias        torch.Size([6])
conv2.weight      torch.Size([16, 6, 5, 5])
conv2.bias        torch.Size([16])
fc1.weight        torch.Size([120, 400])
fc1.bias          torch.Size([120])
fc2.weight        torch.Size([84, 120])
fc2.bias          torch.Size([84])
fc3.weight        torch.Size([10, 84])
fc3.bias          torch.Size([10])
```


Main 2

```
for n_epi in range(10000):
    epsilon = max(0.01, 0.08 - 0.01*(n_epi/200)) #Linear annealing from 8% to 1%
    s = env.reset()

    for t in range(600):
        a = q.sample_action(torch.from_numpy(s).float(), epsilon)
        s_prime, r, done, info = env.step(a)
        done_mask = 0.0 if done else 1.0
        memory.put((s,a,r/200.0,s_prime, done_mask))
        s = s_prime

        if done:
            break
```

```
if memory.size()>2000:
    train(q, q_target, memory, gamma, optimizer, batch_size)

if n_epi%20==0 and n_epi!=0:
    q_target.load_state_dict(q.state_dict())
    print("# of episode :{}, Avg timestep : {:.1f}, buffer size : {}, epsilon : {:.1f}%".format(
        n_epi, avg_t/20.0, memory.size(), epsilon*100))
```

```

48 def train(q, q_target, memory, gamma, optimizer, batch_size):
49     for i in range(10):
50         batch = memory.sample(batch_size)
51         s_lst, a_lst, r_lst, s_prime_lst, done_mask_lst = [], [], [], [], []
52
53         for transition in batch:
54             s, a, r, s_prime, done_mask = transition
55             s_lst.append(s)
56             a_lst.append([a])
57             r_lst.append([r])
58             s_prime_lst.append(s_prime)
59             done_mask_lst.append([done_mask])
60
61         s, a, r, s_prime, done_mask = torch.tensor(s_lst, dtype=torch.float), torch.tensor(a_lst), \
62                                     torch.tensor(r_lst), torch.tensor(s_prime_lst, dtype=torch.float), \
63                                     torch.tensor(done_mask_lst)
64         q_out = q(s)
65         q_a = q_out.gather(1, a)
66         max_q_prime = q_target(s_prime).max(1)[0].unsqueeze(1)
67         target = r + gamma * max_q_prime * done_mask
68         loss = F.smooth_l1_loss(target, q_a)
69
70         optimizer.zero_grad()
71         loss.backward()
72         optimizer.step()

```

Shape : [32,2]

취한 action의 q값만 골라냄.
Shape : [32,1]

$$r + \gamma \max_{a'} Q(s', a'; \theta_i^-)$$

학습 결과

```
# of episode :20, Avg timestep : 9.1, buffer size : 202, epsilon : 7.9%
# of episode :40, Avg timestep : 8.6, buffer size : 393, epsilon : 7.8%
# of episode :60, Avg timestep : 8.7, buffer size : 587, epsilon : 7.7%
# of episode :80, Avg timestep : 8.6, buffer size : 779, epsilon : 7.6%
# of episode :100, Avg timestep : 8.7, buffer size : 973, epsilon : 7.5%
# of episode :120, Avg timestep : 8.6, buffer size : 1165, epsilon : 7.4%
# of episode :140, Avg timestep : 8.8, buffer size : 1360, epsilon : 7.3%
# of episode :160, Avg timestep : 9.0, buffer size : 1560, epsilon : 7.2%
# of episode :180, Avg timestep : 8.7, buffer size : 1754, epsilon : 7.1%
# of episode :200, Avg timestep : 8.7, buffer size : 1947, epsilon : 7.0%
# of episode :220, Avg timestep : 10.6, buffer size : 2179, epsilon : 6.9%
# of episode :240, Avg timestep : 14.6, buffer size : 2491, epsilon : 6.8%
# of episode :260, Avg timestep : 10.9, buffer size : 2729, epsilon : 6.7%
# of episode :280, Avg timestep : 9.9, buffer size : 2947, epsilon : 6.6%
# of episode :300, Avg timestep : 17.4, buffer size : 3316, epsilon : 6.5%
# of episode :320, Avg timestep : 109.5, buffer size : 5525, epsilon : 6.4%
# of episode :340, Avg timestep : 121.2, buffer size : 7970, epsilon : 6.3%
# of episode :360, Avg timestep : 213.4, buffer size : 12259, epsilon : 6.2%
# of episode :380, Avg timestep : 179.5, buffer size : 15869, epsilon : 6.1%
# of episode :400, Avg timestep : 114.1, buffer size : 18171, epsilon : 6.0%
# of episode :420, Avg timestep : 101.2, buffer size : 20215, epsilon : 5.9%
# of episode :440, Avg timestep : 139.5, buffer size : 23025, epsilon : 5.8%
# of episode :460, Avg timestep : 163.5, buffer size : 26315, epsilon : 5.7%
# of episode :480, Avg timestep : 199.3, buffer size : 30322, epsilon : 5.6%
# of episode :500, Avg timestep : 268.6, buffer size : 35715, epsilon : 5.5%
# of episode :520, Avg timestep : 260.0, buffer size : 40935, epsilon : 5.4%
# of episode :540, Avg timestep : 207.1, buffer size : 45096, epsilon : 5.3%
# of episode :560, Avg timestep : 234.8, buffer size : 49812, epsilon : 5.2%
# of episode :580, Avg timestep : 209.8, buffer size : 50000, epsilon : 5.1%
```