

Policy Gradient 구현

reinforce, TD Actor-Critic

[쉽게 구현하는 강화학습 1화]

팡요랩 - 노승은, 전민영

2019.04.27

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

복습 – Policy Gradient (7장)

- The **policy gradient** has many equivalent forms

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \, v_t] && \text{REINFORCE} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \, Q^w(s, a)] && \text{Q Actor-Critic} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \, A^w(s, a)] && \text{Advantage Actor-Critic} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \, \delta] && \text{TD Actor-Critic} \\ &= \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) \, \delta e] && \text{TD}(\lambda) \text{ Actor-Critic} \\ G_{\theta}^{-1} \nabla_{\theta} J(\theta) &= w && \text{Natural Actor-Critic} \end{aligned}$$

- Each leads a stochastic gradient ascent algorithm
- Critic uses **policy evaluation** (e.g. MC or TD learning) to estimate $Q^{\pi}(s, a)$, $A^{\pi}(s, a)$ or $V^{\pi}(s)$

복습 – Policy Gradient (7장)

- For the true value function $V^{\pi_\theta}(s)$, the TD error δ^{π_θ}

$$\delta^{\pi_\theta} = r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$$

- is an unbiased estimate of the advantage function

$$\begin{aligned}\mathbb{E}_{\pi_\theta} [\delta^{\pi_\theta} | s, a] &= \mathbb{E}_{\pi_\theta} [r + \gamma V^{\pi_\theta}(s') | s, a] - V^{\pi_\theta}(s) \\ &= Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s) \\ &= A^{\pi_\theta}(s, a)\end{aligned}$$

- So we can use the TD error to compute the policy gradient

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s, a) \delta^{\pi_\theta}]$$

- In practice we can use an approximate TD error

$$\delta_v = r + \gamma V_v(s') - V_v(s)$$

AutoDiff 이용하기

구해야 하는 것 : $\nabla_{\theta} \log \pi_{\theta}(s, a) r = \mathbf{r} \times \nabla_{\theta} \log \pi_{\theta}(s, a)$

누구를 미분하면 위에 식이 되지? $\mathbf{r} \times \log \pi_{\theta}(s, a)$

그러면 loss 는 어떻게? $-\mathbf{r} \times \log \pi_{\theta}(s, a)$

- 는 왜불지 ? Loss 는 자동으로 minimize 되는데 우리는 maximize 하고 싶으니까!

REINFORCE

```
1 #REINFORCE
2 import gym
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 import torch.optim as optim
7 from torch.distributions import Categorical
```

Import 부터!

REINFORCE

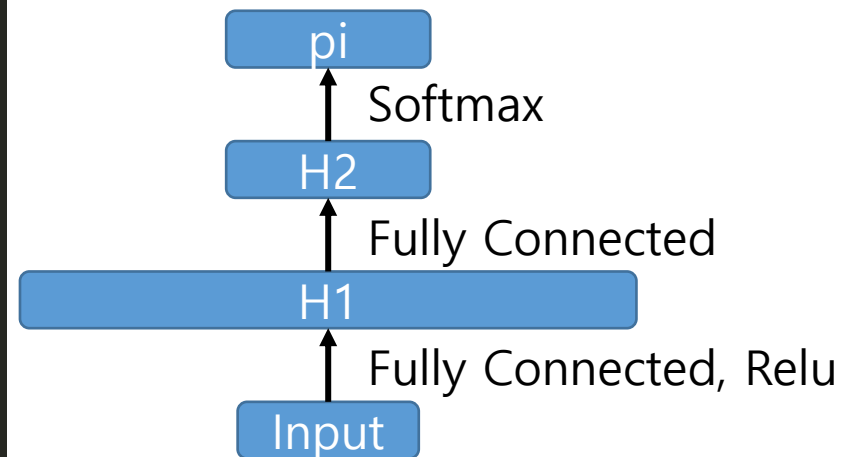
```
37 def main():
38     env = gym.make('CartPole-v1')
39     pi = Policy()
40     avg_t = 0
41
42     for n_epi in range(10000):
43         obs = env.reset()
44         for t in range(600):
45             obs = torch.tensor(obs, dtype=torch.float)
46             out = pi(obs)
47             m = Categorical(out)
48             action = m.sample()
49             obs, r, done, info = env.step(action.item())
50             pi.put_data((r, torch.log(out[action])))
51             if done:
52                 break
53             avg_t += t
54             pi.train()
55             if n_epi%20==0 and n_epi!=0:
56                 print("# of episode : {}, Avg timestep : {}".format(n_epi, avg_t/20.0))
57                 avg_t = 0
58     env.close()
```

Annotations:

- Arrow from `obs` in line 45 points to `tensor([0.0255, -0.0159, -0.0489, -0.0408])`
- Arrow from `out` in line 46 points to `tensor(0)`
- Arrow from `torch.log(out[action])` in line 50 points to `$\log \pi_{\theta}(s, a)$`

REINFORCE

```
9 class Policy(nn.Module):
10     def __init__(self):
11         super(Policy, self).__init__()
12         self.data = []
13         self.gamma = 0.99
14
15         self.fc1 = nn.Linear(4, 128)
16         self.fc2 = nn.Linear(128, 2)
17         self.optimizer = optim.Adam(self.parameters(), lr=0.0005)
18
19     def forward(self, x):
20         x = F.relu(self.fc1(x))
21         x = F.softmax(self.fc2(x), dim=0)
22         return x
23
24     def put_data(self, item):
25         self.data.append(item)
26
27     def train(self):
28         R = 0
29         for r, log_prob in self.data[::-1]:
30             R = r + R * self.gamma
31             loss = -log_prob * R
32             self.optimizer.zero_grad()
33             loss.backward()
34             self.optimizer.step()
35         self.data = []
36
```



$$\log \pi_{\theta}(s, a) V_t$$

TD Actor Critic

```
43 for n_epi in range(10000):
44     obs = env.reset()
45     loss_lst = []
46     for t in range(600):
47         obs = torch.from_numpy(obs).float()
48         pi, v = model(obs)
49         m = Categorical(pi)
50         action = m.sample()
51
52         obs, r, done, info = env.step(action.item())
53         _, next_v = model(torch.from_numpy(obs).float())
54         delta = r + gamma * next_v - v
55         loss = -torch.log(pi[action]) * delta.item() + delta * delta
56         model.gather_loss(loss)
57
58     if done:
59         break
60
61 model.train()
```

$$\delta^{\pi_{\theta}} = r + \gamma V^{\pi_{\theta}}(s') - V^{\pi_{\theta}}(s)$$

$$\log \pi_{\theta}(s, a) \delta^{\pi_{\theta}}$$

Value loss

TD Actor Critic

```
9 class ActorCritic(nn.Module):
10     def __init__(self):
11         super(ActorCritic, self).__init__()
12         self.loss_lst = []
13
14         self.fc1 = nn.Linear(4, 128)
15         self.fc_pi = nn.Linear(128, 2)
16         self.fc_v = nn.Linear(128, 1)
17         self.optimizer = optim.Adam(self.parameters(), lr=0.002)
18
19     def forward(self, x):
20         x = F.relu(self.fc1(x))
21         pol = self.fc_pi(x)
22         pi = F.softmax(pol, dim=0)
23         v = self.fc_v(x)
24         return pi, v
25
26     def gather_loss(self, loss):
27         self.loss_lst.append(loss.unsqueeze(0))
28
29     def train(self):
30         loss = torch.cat(self.loss_lst).sum()
31         loss = loss/len(self.loss_lst)
32         self.optimizer.zero_grad()
33         loss.backward()
34         self.optimizer.step()
35         self.loss_lst = []
36
```

