# Day 3. Policy Gradient and Actor-Critic

**NPEX Reinforcement Learning** 

July 28, 2021 Jaeuk Shin, Minkyu Park



# **Policy Gradient - Recap**

Cast Reinforcement learning problem into optimization problem

$$\max_{\theta} J(\theta) := E_{\tau \sim p_{\theta}(\tau)} \left[ \sum \gamma^t r_t \right]$$

Where policy is represented with parameter  $\theta \rightarrow \pi_{\theta}$ 

We can solve above optimization problem with gradient ascent

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$
How?



Recall the gradient of parametric policy  $\pi_{\theta}$  for objective function J

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right) \right]$$



Recall the gradient of parametric policy  $\pi_{\theta}$  for objective function J

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right) \right]$$

Approximate with sample mean by sampling N trajectories  $\tau = (s_0, a_0, \dots, s_T, a_T)$ 

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right)$$



Recall the gradient of parametric policy  $\pi_{\theta}$  for objective function J

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right) \right]$$

Approximate with sample mean by sampling N trajectories  $\tau = (s_0, a_0, \dots, s_T, a_T)$ 

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right)$$



REINFORCE Algorithm



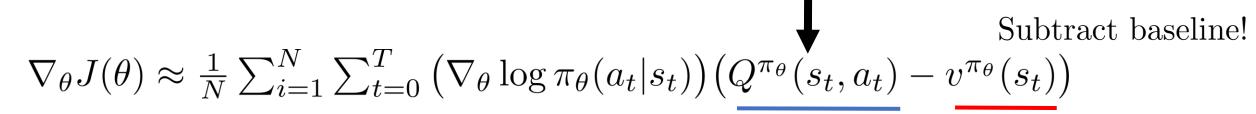
Additional approximation and baseline

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=0}^{T} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \underbrace{\left( \sum_{t=0}^{T} r(s_t, a_t) \right)}_{\text{Const}}$$



$$\nabla J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \left[ \left( \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t=0}^{T} r(s_t, a_t) \right) \right]$$

$$\nabla J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \left[ \left( \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( \sum_{t'=t}^{T} r(s_{t'}, a_{t'}) \right) \right]$$



 $=A^{\pi_{\theta}}(s_t,a_t)$ 



How to get  $\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ ?

```
class Policy(nn.Module):
         def __init__(self, state_dim, num_action, hidden_size1, hidden_size2):
             super(Policy, self). init ()
             self.fc1 = nn.Linear(state dim, hidden size1)
             self.fc2 = nn.Linear(hidden_size1, hidden_size2)
                                                                      NN returns |A| dimension vector
             self.fc3 = nn.Linear(hidden size2, num action)
        def forward(self, x):
            x = self.fc1(x)
                                        Return softmax value of length |\mathcal{A}|, which can be used as probability to
            x = F.relu(x)
10
                                    choose action \to \pi_{\theta}
            x = self.fc2(x)
11
            x = F.relu(x)
12
             action score = self.fc3(x)
13
             return F.softmax(action score, dim=1)
14
15
```



How to get  $\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ ?

```
def select action(state):
        state = torch.from numpy(state).float().unsqueeze(0)
        state = state.to(device)
                                    → NN output - probability for each action
        probs = pi(state)
       m = Categorical(probs)
6
        action = m\sample()
        return action.item(), m.log_prob(action)
9
```

This function builds probability distribution for discrete action space and give  $\log \pi_{\theta}(a_t|s_t)$ ,  $\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$ 



#### Create $\nabla_{\theta} J(\theta)$

return pi loss / num trajs

21

```
def calculate_PG(pi_returns_discounted, dataset):
         pi loss = 0
         for data in dataset:
                                          Reversed order
              advantage, DCR = [], 0
              for r in reversed(data['reward']):
                  # TODO : Caculate discounted redataset.append(data)turn from t=i
                  # Hint : reversed() will give saved rewards in reversed order
                                                                                       Q^{\pi_{\theta}}(s_t, a_t) \approx \sum_{t=t'}^{T} \gamma^{t-t'} r(s_{t'}, a_{t'})
                  DCR = r + gamma * DCR
                  \# Q(s,a) is replaced with discounted sum of rewards (DCR)
10
                  \# v(s) is replaced with empirical v(s 0)
                                                                                               \rightarrow baseline with E_{s_0,\pi_{\theta}} \left[ \sum_{t=0}^{T} \gamma^t r_t \right]
                  advantage.insert(0, DCR - np.mean(pi returns discounted)
13
              # TODO: alternate between two losses to see difference!
14
              pi loss vanilla = [log pi * DCR for log pi in data['log pi']]
              pi_loss_baseline = [log_pi * a for log_pi, a in zip(data['log pi'], advantage)]
16
17
              # Take mean value
                                                                Choose between loses!
              pi loss += torch.cat(pi loss vanilla).sum()
19
20
```



Using baseline increaes performance dramatically!

```
Epoch 0 Return mean: 22.74
                                Return std: 12.03
                                                         Time
                                                                        Epoch 0 Return mean: 19.76
                                                                                                         Return std: 8.49
                                                                                                                                 Time(I
                                Return std: 13.00
                                                         Time
Epoch 5 Return mean: 25.18
                                                                        Epoch 5 Return mean: 29.37
                                                                                                         Return std: 18.32
                                                                                                                                 Time(
Epoch 10
                Return mean: 30.49
                                         Return std: 15.26
                                                                        Epoch 10
                                                                                        Return mean: 40.56
                                                                                                                 Return std: 25.14
Epoch 15
                                         Return std: 19.31
                Return mean: 38.63
                                                                        Epoch 15
                                                                                        Return mean: 53.07
                                                                                                                 Return std: 26.29
Epoch 20
                Return mean: 44.66
                                         Return std: 22.34
                                                                                        Return mean: 76.31
                                                                        Epoch 20
                                                                                                                 Return std: 44.41
Epoch 25
                Return mean: 60.22
                                         Return std: 30.04
                                                                        Epoch 25
                                                                                        Return mean: 131.40
                                                                                                                 Return std: 71.74
Epoch 30
                Return mean: 57.98
                                         Return std: 28.14
                                                                        Epoch 30
                                                                                        Return mean: 192.59
                                                                                                                 Return std: 94.73
Epoch 35
                Return mean: 66.45
                                         Return std: 31.69
                                                                        Epoch 35
                                                                                        Return mean: 279.06
                                                                                                                 Return std: 117.43
                                         Return std: 43.28
Epoch 40
                Return mean: 85.14
                                                                                        Return mean: 325.63
                                                                                                                 Return std: 113.61
                                                                        Epoch 40
Epoch 45
                Return mean: 111.59
                                         Return std: 59.01
                                                                                        Return mean: 399.72
                                                                        Epoch 45
                                                                                                                 Return std: 122.64
Epoch 50
                Return mean: 126.32
                                         Return std: 55.62
                                                                                        Return mean: 429.90
                                                                        Epoch 50
                                                                                                                 Return std: 97.34
Epoch 55
                Return mean: 155.78
                                         Return std: 65.45
                                                                        Epoch 55
                                                                                        Return mean: 441.48
                                                                                                                 Return std: 86.38
Epoch 60
                Return mean: 191.32
                                         Return std: 77.60
                                                                        Epoch 60
                                                                                        Return mean: 456.22
                                                                                                                 Return std: 86.65
                Return mean: 205.11
Epoch 65
                                         Return std: 79.55
                                                                                        Return mean: 474.78
                                                                        Epoch 65
                                                                                                                 Return std: 66.90
Epoch 70
                Return mean: 208.84
                                         Return std: 88.27
                                                                                        Return mean: 493.40
                                                                        Epoch 70
                                                                                                                 Return std: 33.79
Epoch 75
                Return mean: 158.31
                                         Return std: 53.31
                                                                                        Return mean: 483.57
                                                                        Epoch 75
                                                                                                                 Return std: 57.95
Epoch 80
                Return mean: 141.18
                                         Return std: 53.16
                                                                        Epoch 80
                                                                                        Return mean: 489.79
                                                                                                                 Return std: 44.05
Epoch 85
                Return mean: 193.22
                                         Return std: 58.09
                                                                                                     479.35
                                                                        Epoch 85
                                                                                        Return mean:
                                                                                                                 Return std: 59.41
Epoch 90
                Return mean: 160.01
                                         Return std: 45.93
                                                                        Epoch 90
                                                                                        Return mean: 488.63
                                                                                                                 Return std: 41.79
Epoch 95
                Return mean: 307.46
                                         Return std: 101.98
                                                                        Epoch 95
                                                                                        Return mean: 493.34
                                                                                                                 Return std: 34.31
```

Without baseline

Without empirical mean baseline



Actor-Critic architecture

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \left( \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right) \left( Q^{\pi_{\theta}}(s_t, a_t) - v^{\pi_{\theta}}(s_t) \right)$$

Fit advantage  $A^{\pi}(s, a) = Q^{\pi}(s, a) - v^{\pi}(s)$  with parameter w!

With value function estimation  $v^{\pi} \approx v_w^{\pi}$ ,  $\to A^{\pi}(s, a) = r + \gamma v_w^{\pi}(s') - v_w^{\pi}(s)$ 

Here, w is called critic,  $\theta$  is called actor



Example: Linear critic

Define linear feature set as follow for CartPole task

$$\phi(s) = (1, e^{-\frac{\|s-\mu_1\|^2}{2\sigma^2}}, \cdots, e^{-\frac{\|s-\mu_M\|^2}{2\sigma^2}})$$
 (radial basis functions)

 $v_w^{\pi}(s) = \phi(s)^{\top} w$ , where w is parameter vector of length M+1

Perform linear fitting with sampled (s, a, r, s')

$$x = [\phi(s^0), \cdots, \phi(s^N)]$$

$$y = [r^0 + \gamma \phi(s'^0)^\top w_{old}, \cdots, r^N + \gamma \phi(s'^N)^\top w_{old}]$$

Fit  $w_{new}$  to minimize  $||w^{\top}x - y||_2^2$ 



#### Example: Linear critic

```
def calculate vf(dataset, vf):
         X, y = [], []
          for data in dataset:
              for s, next_s, r in zip(data['state'], data['next_state'], data['reward']):
                   v = state2feature(s)
                                                            \phi(s')^{\top}w_{old}
                  0 = r
                   if vf is not None:
                       Q = r + gamma * vf.predict(state2feature(next_s).reshape(1, -1))[0]
                   X.append(v)
10
                                          x = [\phi(s^0), \cdots, \phi(s^N)]
11
                   y.append(Q)
12
                                           y = [r^0 + \gamma \phi(s'^0)^\top w_{old}, \cdots, r^N + \gamma \phi(s'^N)^\top w_{old}]
          return X, y
13
```



#### Algorithm summary

Sample trajectory  $(s_0, a_0, \dots, s_T, a_T)$ 

Fit w with TD error

Perform policy gradient with 
$$(\nabla_{\theta} \log \pi_{\theta}(a_t|s_t))(r + \gamma v_w^{\pi}(s_{t+1}) - v_w^{\pi}(s_t))$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=0}^{T} \left( \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \right) \left( Q^{\pi_{\theta}}(s_t, a_t) - v^{\pi_{\theta}}(s_t) \right)$$

Note, we don't need full trajectory now (incremental online algorithm)



#### Algorithm summary

```
def get advantage(data, vf):
          advantage, baseline = [], [] \phi(s)^\top w \qquad \qquad \phi(s) for s, next_s, r in zip(data['state'], data['next_state'], data['reward']): v = vf.predict(state2feature(s).reshape(1, -1))[0]
                 v next = vf.predict(state2feature(next s).reshape(1, -1))[0]
                 # TODO: Complete advantage calculation by calculating Q-value
                Q = r + gamma * v_next
                 (r + \gamma v_w^{\pi}(s_{t+1}) - v_w^{\pi}(s_t)) 
10
                 advantage.append(A)
11
                 baseline.append(v)
12
13
           return advantage, baseline
14
```



Slower learning due to linear fitting calculation time

```
Epoch 0 Return mean: 19.53
                                                                                                        Return std: 8.33
                                                                                                                                 Time
Epoch 0 Return mean: 19.76
                                Return std: 8.49
                                                         Time(I
                                                                       Epoch 5 Return mean: 29.72
                                                                                                        Return std: 18.20
                                                         Time(
                                                                                                                                 Time
Epoch 5 Return mean: 29.37
                                Return std: 18.32
                                                                       Epoch 10
                                                                                       Return mean: 38.70
                                                                                                                Return std: 24.50
Epoch 10
                Return mean: 40.56
                                        Return std: 25.14
                                                                       Epoch 15
                                                                                       Return mean: 54.86
                                                                                                                Return std: 32.84
Epoch 15
                Return mean: 53.07
                                        Return std: 26.29
                                                                       Epoch 20
                                                                                       Return mean: 72.85
                                                                                                                Return std: 34.04
                                        Return std: 44.41
Epoch 20
                Return mean: 76.31
                                                                       Epoch 25
                                                                                        Return mean: 121.04
                                                                                                                Return std: 67.89
Epoch 25
                Return mean: 131.40
                                        Return std: 71.74
                                                                       Epoch 30
                                                                                       Return mean: 206.30
                                                                                                                Return std: 101.12
Epoch 30
                Return mean: 192.59
                                        Return std: 94.73
                                                                       Epoch 35
                                                                                       Return mean: 268.13
                                                                                                                Return std: 112.62
Epoch 35
                Return mean: 279.06
                                        Return std: 117.43
                                                                                        Return mean: 343.16
                                                                       Epoch 40
                                                                                                                Return std: 110.35
Epoch 40
                Return mean: 325.63
                                        Return std: 113.61
                                                                       Epoch 45
                                                                                        Return mean: 411.13
                                                                                                                Return std: 95.30
                Return mean: 399.72
                                        Return std: 122.64
Epoch 45
                                                                                        Return mean: 428.35
                                                                       Epoch 50
                                                                                                                Return std: 100.25
Epoch 50
                Return mean: 429.90
                                        Return std: 97.34
                                                                                        Return mean: 446.02
                                                                       Epoch 55
                                                                                                                Return std: 91.36
Epoch 55
                Return mean: 441.48
                                        Return std: 86.38
                                                                                        Return mean: 476.90
                                                                                                                Return std: 57.94
                                                                       Epoch 60
                                        Return std: 86.65
                Return mean: 456.22
Epoch 60
                                                                       Epoch 65
                                                                                        Return mean: 477.56
                                                                                                                Return std: 55.69
                Return mean: 474.78
                                        Return std: 66.90
Epoch 65
                                                                       Epoch 70
                                                                                        Return mean: 495.10
                                                                                                                Return std: 21.86
Epoch 70
                Return mean: 493.40
                                        Return std: 33.79
                                                                                        Return mean: 468.04
                                                                                                                Return std: 64.06
                                                                       Epoch 75
Epoch 75
                Return mean: 483.57
                                        Return std: 57.95
                                                                                        Return mean: 483.19
                                                                       Epoch 80
                                                                                                                Return std: 50.66
                             489.79
                                        Return std: 44.05
Epoch 80
                Return mean:
                                                                       Epoch 85
                                                                                        Return mean: 475.75
                                                                                                                Return std: 62.50
                Return mean: 479.35
Epoch 85
                                        Return std: 59.41
                                                                       Epoch 90
                                                                                        Return mean: 492.98
                                                                                                                Return std: 28.44
Epoch 90
                Return mean: 488.63
                                        Return std: 41.79
                                                                                                                Return std: 19.28
                                                                       Epoch 95
                                                                                        Return mean: 497.29
                Return mean: 493.34
                                        Return std: 34.31
Epoch 95
```

With average reward baseline

With linear fitted value function baseline

