

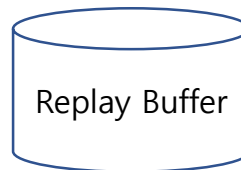
# 최적화 학습 알고리즘

## 코드 수정

2023/06/01

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.



# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$

[/run.py#L418~457](#)에 해당

2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$

3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$

4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$

5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$

1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$

2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$

3.  $\theta_{t+1} = \theta^K$

6. Repeat.

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/run.py#L418](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/run.py#L418)

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$

2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$

3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$

4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$

5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$

1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$

2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$

3.  $\theta_{t+1} = \theta^K$

6. Repeat.

ppoTF2.replayNew(): /ppoTF2.py#L765~805에 해당

run.py#L481 → ppoTF2.py#L690 → ppoTF2.py#L765 순으로 호출.

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/run.py#L481](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/run.py#L481)

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L690](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L690)

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L765](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L765)

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Target value 계산: /ppoTF2.py#L790 and #795

Training: /ppoTF2.py# 804

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L790](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L790)  
[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L795](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L795)  
[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L804](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L804)

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Critic DNN Model: /ppoTF2.py#249

Critic Loss: /ppoTF2.py#282

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L249](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L249)  
[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L282](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L282)

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}', \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i')$

4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i') - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$

Advantage 계산: /ppoTF2.py#795

5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L795](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L795)

# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Training: /ppoTF2.py#803

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent\\_tf2/policy/ppoTF2.py#L803](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atasc-rl/multiagent_tf2/policy/ppoTF2.py#L803)



# 현재 알고리즘

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$  using target  $y_i = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = r_i + \gamma \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i | \mathbf{s}_i)}{\pi(\mathbf{a}_i | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Actor DNN Model: /ppoTF2.py#180

PPO Loss  $J_t(\theta)$ : /ppoTF2.py#217

[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atssc-rl/multiagent\\_tf2/policy/ppoTF2.py#L180](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atssc-rl/multiagent_tf2/policy/ppoTF2.py#L180)  
[https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atssc-rl/multiagent\\_tf2/policy/ppoTF2.py#L217](https://github.com/etri-city-traffic-brain/traffic-signal-optimization/blob/8ac45baeedda25f78d5343223ac2edee7aa44ee8/atssc-rl/multiagent_tf2/policy/ppoTF2.py#L217)

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

이전 코드와 동일. /run.py#L418~457에 해당

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

이전 코드와 동일. /ppo.py#L1032~1041에 해당

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Critic DNN Model: </ppo.py#L328>

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_{\theta}(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Value for next state  $\mathbf{s}'_i$

$$\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i) \approx \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i)$$

$\mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$

[/ppo.py#L1060, #959](#)

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Compute target values. /ppo.py#1066

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_{\theta}(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Training: [/ppo.py#1072](#)



# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i)$ ;  $\mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$

4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t)$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$

[/ppo.py#978](#)

5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$

1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$

2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$

3.  $\theta_{t+1} = \theta^K$

6. Repeat.

$$\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i)$$

$$\hat{V}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i) \approx \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t) \leftarrow \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$$

Sample actions according to the current policy.

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i)$ ;  $\mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t)$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Actor DNN Model: /ppoTF2.py#172

PPO Loss  $J_t(\theta)$ : /ppoTF2.py#233

기준과 동일

# 수정 알고리즘(Off-Policy)

Actually, we do not need this one.

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i)$ ;  $\mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t)$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}$ ;  $\mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

/ppo.py#1077 and #995

Augment a batch by sampling actions according to the current policy.

# 수정 알고리즘(Off-Policy)

1. Take action  $\mathbf{a} \sim \pi(\cdot | \mathbf{s})$ , get  $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}', \log \pi(\mathbf{a} | \mathbf{s}))$  and store in  $R$
2. Sample a batch  $\{\mathbf{s}_i, \mathbf{a}_i, r_i, \mathbf{s}'_i, \log \pi(\mathbf{a}_i | \mathbf{s}_i)\}$  from buffer  $R$
3. Update  $\hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i)$  using target  $y_i = r_i + \gamma \frac{1}{M} \sum_{\mathbf{a}'_i} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}'_i, \mathbf{a}'_i); \mathbf{a}'_i \sim \pi_{\theta_t}(\cdot | \mathbf{s}'_i)$
4. Evaluate  $\hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) = \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i) - \frac{1}{M} \sum_{\mathbf{a}_i^t} \hat{Q}_\phi^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t); \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$
5. Minibatch Learning on  $\{\mathbf{s}_i, \mathbf{a}_i^t, \log \pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)\}; \mathbf{a}_i^t \sim \pi_{\theta_t}(\cdot | \mathbf{s}_i)$ 
  1.  $J_t(\theta) = \frac{1}{N} \sum_i \min \left( \frac{\pi_\theta(\mathbf{a}_i^t | \mathbf{s}_i)}{\pi_{\theta_t}(\mathbf{a}_i^t | \mathbf{s}_i)} \hat{A}^{\pi_{\theta_t}}(\mathbf{s}_i, \mathbf{a}_i^t), \sim \right)$  and  $\theta^1 = \theta_t$
  2. for  $k = 1, \dots, K$  do  $\theta^{k+1} \leftarrow \theta^k + \alpha \nabla_{\theta^k} J_t(\theta^k)$
  3.  $\theta_{t+1} = \theta^K$
6. Repeat.

Training: /ppo.py#1089

Q & A