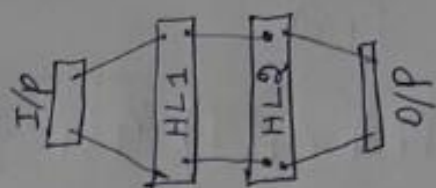


# Policy Gradient



Policy Net - Q Network - NN

→ Q-learning → Minimizing the F1 smooth loss. means, it is working for maximizing two consecutive action cumulative reward.

→ Policy Gradient → Sampling the episodes and optimizing the policy net by gradient ascent.

Objective Function:  $J(\theta) = E \left[ \sum_{t=0}^{T-1} r_{t+1} \right] \rightarrow ①$

here, nomenclature:  $\{S_t, a_t, S_{t+1}, r_{t+1}\}$

Gradient Ascent:  $\theta \leftarrow \theta + \frac{\partial}{\partial \theta} [J(\theta)] \dots \dots ②$

As we all know that,  $E[f(x)] = \sum_x P(x) f(x)$ .

$$J(\theta) = E \left[ \sum_{t=0}^{T-1} r_{t+1} \mid \pi_{\theta} \right]$$

$$J(\theta) = \sum_{t=0}^{T-1} P(S_t, a_t \mid \tau) r_{t+1} \dots \dots \dots ③$$

i → Arbitrary Starting Point

$\tau$  → Given Trajectory

Differentiating both sides with respect to policy parameter  $\theta$ ,  
 using  $\frac{d}{dx} [\log f(x)] = \frac{f'(x)}{f(x)}$  ----- (4)

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \sum_{t=0}^{T-1} \nabla_{\theta} P(s_t, a_t | \tau) r_{t+1} \\ &= \sum_{t=0}^{T-1} P(s_t, a_t | \tau) \cdot \frac{\nabla_{\theta} P(s_t, a_t | \tau)}{P(s_t, a_t | \tau)} r_{t+1} \\ &= \sum_{t=0}^{T-1} P(s_t, a_t | \tau) \cdot \nabla_{\theta} \log P(s_t, a_t | \tau) \cdot r_{t+1}\end{aligned}$$

So, substituting in the equation (4),

$$\nabla_{\theta} J(\theta) = E \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log P(s_t, a_t | \tau) r_{t+1} \right]$$

This approximate can be done by re-writing the equation,

$$\nabla_{\theta} J(\theta) \sim \sum_{t=0}^{T-1} \nabla_{\theta} \log P(s_t, a_t | \tau) r_{t+1} \text{ ----- (5)}$$

looking expression for  $\nabla_{\theta} \log P(s_t, a_t | \tau)$ ,

$$\begin{aligned}P(s_t, a_t | \tau) &= P(s_0, a_0, s_1, a_1, \dots, s_{t-1}, a_{t-1}, s_t, a_t | \pi_{\theta}) \\ &= P(s_0) \pi_{\theta}(a_1 | s_0) P(s_1 | s_0, a_0) \pi_{\theta}(a_2 | s_1) P(s_2 | s_1, a_1) \\ &\quad \pi_{\theta}(a_3 | s_2) P(s_3 | s_2, a_2) \dots \dots \dots \\ &\quad \dots \dots P(s_{t-1} | s_{t-2}, a_{t-2}) \pi_{\theta}(a_{t-1} | s_{t-2}) \\ &\quad \dots \dots \dots P(s_t | s_{t-1}, a_{t-1}) \pi_{\theta}(a_t | s_{t-1})\end{aligned}$$

here, the function term,  $\nabla_{\theta} \log P(s_t, a_t | \tau)$ ,

$$\begin{aligned} \nabla_{\theta} \log P(s_t, a_t | \tau) &= \nabla_{\theta} \log P(s_0) + \nabla_{\theta} \log \pi_{\theta}(a_0 | s_0) + \\ &\quad \nabla_{\theta} \log P(s_1 | s_0, a_0) + \nabla_{\theta} \log \pi_{\theta}(a_1 | s_1) + \dots \\ &\quad \dots + \nabla_{\theta} \log P(s_{t-1} | s_{t-2}, a_{t-2}) \\ &\quad + \nabla_{\theta} \log \pi_{\theta}(a_{t-1} | s_{t-1}) + \dots \\ &\quad \dots + \nabla_{\theta} \log P(s_t, s_{t-1}, a_{t-1}) + \nabla_{\theta} \log \pi_{\theta}(a_t | s_{t-1}) \end{aligned}$$

Here, it is important to note that  $P(s_t | s_{t-1}, a_{t-1})$  is not dependant on the policy parameter  $\theta$  and solely dependant on environment reinforcement learning.

$$\nabla_{\theta} \log P(s_t, a_t | \tau) = \sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'} | s_{t'})$$

$$\begin{aligned} \text{from, } \nabla_{\theta} \log P(s_t, a_t | \tau) &= 0 + \nabla_{\theta} \log \pi_{\theta}(a_1 | s_0) + 0 + \\ &\quad \nabla_{\theta} \log \pi_{\theta}(a_2 | s_1) + 0 + \dots + \dots \\ &\quad + 0 + \nabla_{\theta} \log \pi_{\theta}(a_{t-1} | s_{t-2}) + \dots \\ &\quad \dots + \nabla_{\theta} \log \pi_{\theta}(a_t | s_{t-1}) \end{aligned}$$

So,

$$\nabla_{\theta} \log P(s_t, a_t | \tau) = \sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'} | s_{t'}) \quad \text{--- (6)}$$

from equation (5) & (6),

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^{T-1} r_{t+1} \left\{ \sum_{t'=0}^t \nabla_{\theta} \log \pi_{\theta}(a_{t'} | s_{t'}) \right\} \quad \text{--- (7)}$$



$$\begin{aligned}
&= r_1 \left[ \sum_{t'=0}^0 \nabla_{\theta} \log \pi_{\theta}(a_{t'} | S_{t'}) \right] + r_2 \left[ \sum_{t'=0}^1 \nabla_{\theta} \log \pi_{\theta}(a_{t'} | S_{t'}) \right] \\
&+ r_3 \left[ \sum_{t'=0}^2 \nabla_{\theta} \log \pi_{\theta}(a_{t'} | S_{t'}) \right] + \dots \\
&\dots + r_{T-1} \left[ \sum_{t'=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{t'} | S_{t'}) \right] \\
&= r_1 \nabla_{\theta} \log \pi_{\theta}(a_0 | S_0) + r_2 \left[ \nabla_{\theta} \log \pi_{\theta}(a_0 | S_0) + \nabla_{\theta} \log \pi_{\theta}(a_1 | S_1) \right] \\
&+ r_3 \left[ \nabla_{\theta} \log \pi_{\theta}(a_0 | S_0) + \nabla_{\theta} \log \pi_{\theta}(a_1 | S_1) + \nabla_{\theta} \log \pi_{\theta}(a_2 | S_2) \right] \\
&\dots + \nabla_{\theta} \log \pi_{\theta}(a_{T-1} | S_{T-1}) \cdot r_T
\end{aligned}$$

$$\begin{aligned}
&= \nabla_{\theta} \log \pi_{\theta}(a_0 | S_0) [r_1 + r_2 + r_3 + \dots + r_T] \\
&+ \nabla_{\theta} \log \pi_{\theta}(a_1 | S_1) [r_2 + r_3 + \dots + r_T] \\
&+ \nabla_{\theta} \log \pi_{\theta}(a_2 | S_2) [r_3 + r_4 + \dots + r_T] + \dots \\
&\dots + \nabla_{\theta} \log \pi_{\theta}(a_{T-1} | S_{T-1}) r_T \\
&\nabla_{\theta} J(\theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | S_t) \left[ \sum_{t'=t+1}^T r_{t'} \right] \text{-----} \textcircled{8}
\end{aligned}$$

Simplifying the term  $\sum_{t'=t+1}^T r_{t'}$  to  $G_t$ .

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | S_t) G_t \text{-----} \textcircled{9}$$

Incorporating the discount factor  $\gamma \in [0, 1]$  into our objective.

$$J(\theta) = E \left[ \gamma^0 r_1 + \gamma^1 r_2 + \gamma^2 r_3 + \dots + \gamma^{T-1} r_{T-1} \mid \pi_{\theta} \right] \text{-----} \textcircled{10}$$

We can perform a similar derivation to obtain,

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left[ \sum_{t'=t+1}^T \gamma^{t'-t-1} r_{t'} \right]$$

and simplifying  $\sum_{t'=t+1}^T \gamma^{t'-t-1} r_{t'}$  to  $G_t$

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t$$

Policy update equation

Pseudo Code:

```
Initialise  $\theta$  arbitrarily
for each episode  $\{s_1, a_1, r_1, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_{\theta}$  do
  for  $t = 1$  to  $T-1$  do
     $\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) V_t$ 
  end for
end for
return  $\theta$ 
```