



Difference between On-policy, Off-policy methods ; Value gradients vs Policy gradients

2021 Summer Internship Program

202011875 김현우





Flow

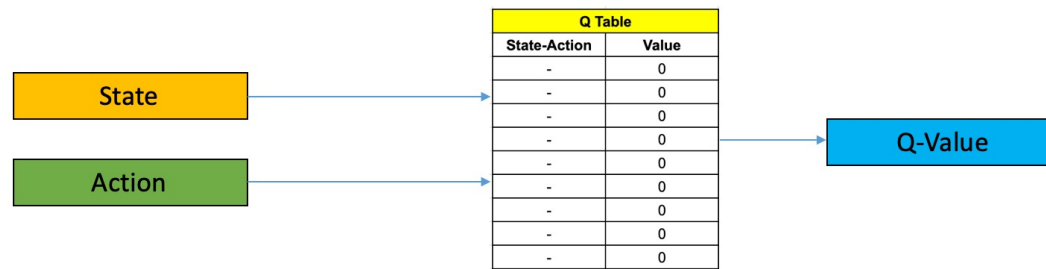
1. *Understanding RL Structure*
2. *Value Gradient advanced examples*
3. Policy Gradient advanced examples



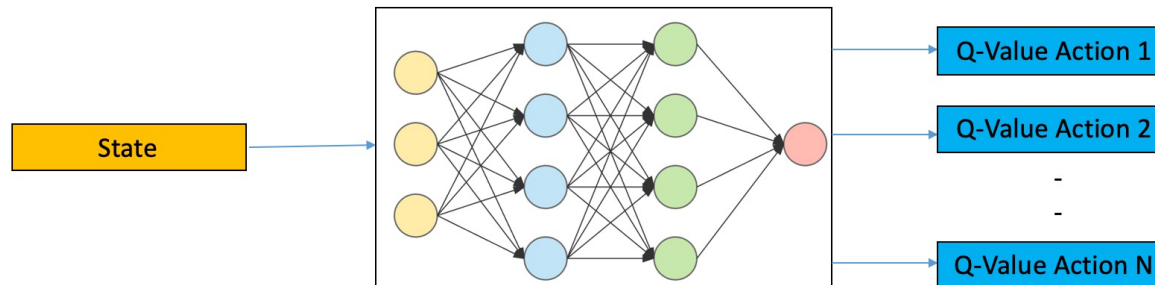
Contents

1. DQN
2. Policy Gradient ; REINFORCE, Actor-Critic
3. CartPole-v1 ; Value gradient vs Policy gradient
4. Pendulum-v0 ; On-policy vs Off-policy
5. Further Study

DQN



Q Learning



Deep Q Learning

DQN

Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory \mathcal{D} to capacity N

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \arg\max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

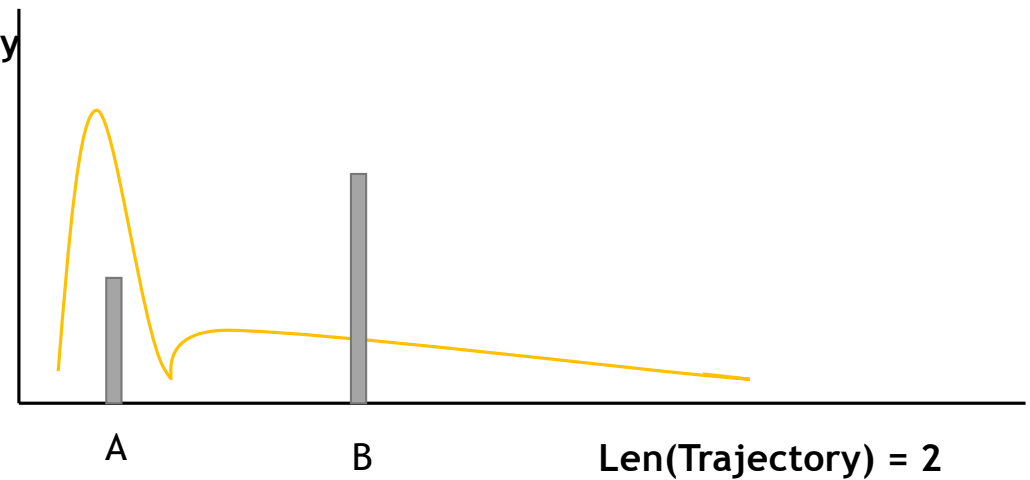
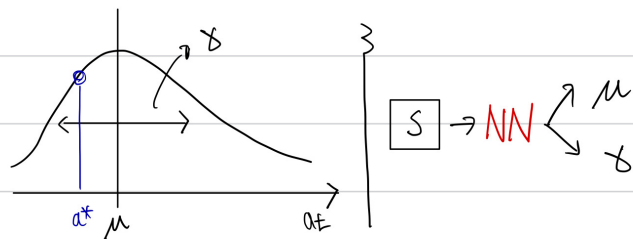
 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

Policy Gradient

$p(a_t|h_t)$ (Value Gradient $\rightarrow \epsilon$ -greedy Probability
Policy Gradient \rightarrow sampling Reward



REINFORCE

$$① J_{\theta} = E[G_t] = \int_t G_t \cdot P_{\theta}(t) dt$$

Find policy maximizes J_{θ} ... Gradient ascent ↗ expectation form

$$② \nabla_{\theta} J_{\theta} = \nabla_{\theta} \int_t (G_t \cdot P_{\theta}(t)) dt = \int_t (G_t \cdot \nabla_{\theta} P_{\theta}(t)) dt = \int_t (G_t \cdot P_{\theta}(t) \cdot \nabla_{\theta} \ln P_{\theta}(t)) dt$$

$$\nabla_{\theta} P_{\theta}(t) = P_{\theta}(t) \cdot \nabla_{\theta} \ln P_{\theta}(t) \dots \nabla_{\theta} \ln P_{\theta}(t) = \frac{\nabla_{\theta} P_{\theta}(t)}{P_{\theta}(t)}$$

$$③ \nabla_{\theta} J_{\theta} = \int_t (G_t \cdot P_{\theta}(t) \cdot \nabla_{\theta} \ln P_{\theta}(a_t | s_t)) dt$$

$$P_{\theta}(t) = P(s_0) \times P(a_0 | s_0) \times P(s_1 | s_0, a_0) \times P(a_1 | s_1) \times \dots \text{Bayesian, Log rule}$$

$$④ \nabla_{\theta} J_{\theta} = \int_t \sum_{k=0}^{\infty} \left[\nabla_{\theta} \ln P_{\theta}(a_t | s_t) \times \sum_{k=t}^{\infty} \gamma^k \cdot R_k \right] P_{\theta}(t) dt$$

$$\begin{aligned} G_t \times \nabla_{\theta} \ln P_{\theta}(a_t | s_t) &= \left[R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \right] \times \nabla_{\theta} \left[\ln P_{\theta}(a_0 | s_0) + \ln P_{\theta}(a_1 | s_1) + \dots \right] \\ \rightarrow \int_t \left[P_{\theta}(t) \times \left(R_0 \times \ln P_{\theta}(a_1 | s_1) \right) + \dots \right] dt &= \int_{t=a_1} R_0 \int_{a_1} \nabla_{\theta} \ln P_{\theta}(a_1 | s_1) da_1 P(t=a_1) dt - a_1 \end{aligned}$$

$$\Rightarrow \nabla_{\theta} \cdot \int_{a_1} P_{\theta}(a_1 | s_1) da_1 = \nabla_{\theta} \cdot 1 = 0$$

$$⑤ \nabla_{\theta} J_{\theta} = \int_t \sum_{k=0}^{\infty} \left(\nabla_{\theta} \ln P_{\theta}(a_t | s_t) \cdot G_t \right) P_{\theta}(t) dt$$

$$\int x \cdot p(x) dx \approx \frac{1}{N} \cdot \sum_{i=1}^N x_i \quad (\text{Sample mean})$$

(N=1) → REINFORCE

$$\textcircled{1} \quad \nabla_{\theta} J_{\theta} = \int_t \sum_{t=0}^{\infty} \nabla_{\theta} \ln P_{\theta}(a_t | h_t) G_t P_{\theta}(t) dt$$

$$P_{\theta}(t) = P_{\theta}(h_{t+1}, a_{t+1}, \dots | h_t, a_t) \times P(h_0, a_0, \dots, h_t, a_t)$$

Actor-Critic

$$\textcircled{2} \quad \nabla_{\theta} J_{\theta} = \int_t \sum_{t=0}^{\infty} \nabla_{\theta} \ln P_{\theta}(a_t | h_t) \int G_t P_{\theta}(h_{t+1}, a_{t+1}, \dots | h_t, a_t) d h_{t+1}, \dots \times P_{\theta}(s_0, \dots, s_t, a_t) d s_0, \dots, a_t$$

$$Q(h_t, a_t) = \int_{s_{t+1}=a_{t+1}} G_t \cdot P_{\theta}(h_{t+1}, a_{t+1}, \dots | h_t, a_t) d h_{t+1}, \dots$$

$$\textcircled{3} \quad \nabla_{\theta} J_{\theta} = \sum_{t=0}^{\infty} \int_{s_t, a_t} \nabla_{\theta} \ln P_{\theta}(a_t | h_t) Q(h_t, a_t) P_{\theta}(a_t | s_t) d h_t, a_t$$

\Rightarrow Return \sim "Q" : Episodic Free \Rightarrow TD 학습

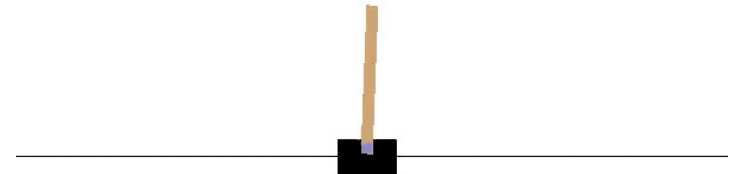
$$\left[\begin{array}{l} \theta \leftarrow \theta + \nabla_{\theta} \ln P_{\theta}(a_t | h_t) Q(h_t, a_t) \\ w \leftarrow w - \beta \nabla_w (R_0 + \gamma Q_w(h_{t+1}, a_{t+1}) - Q_w(h_t, a_t))^2 \end{array} \right.$$

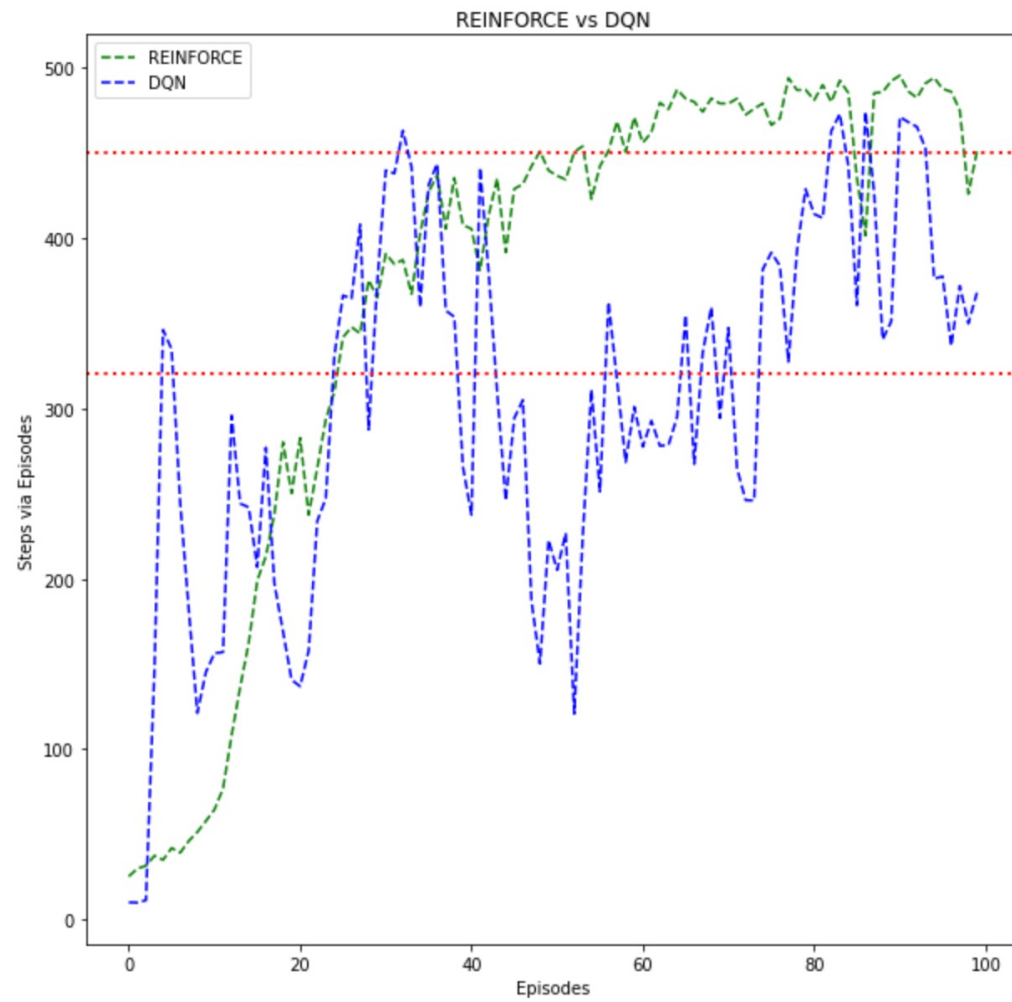
CartPole-v1

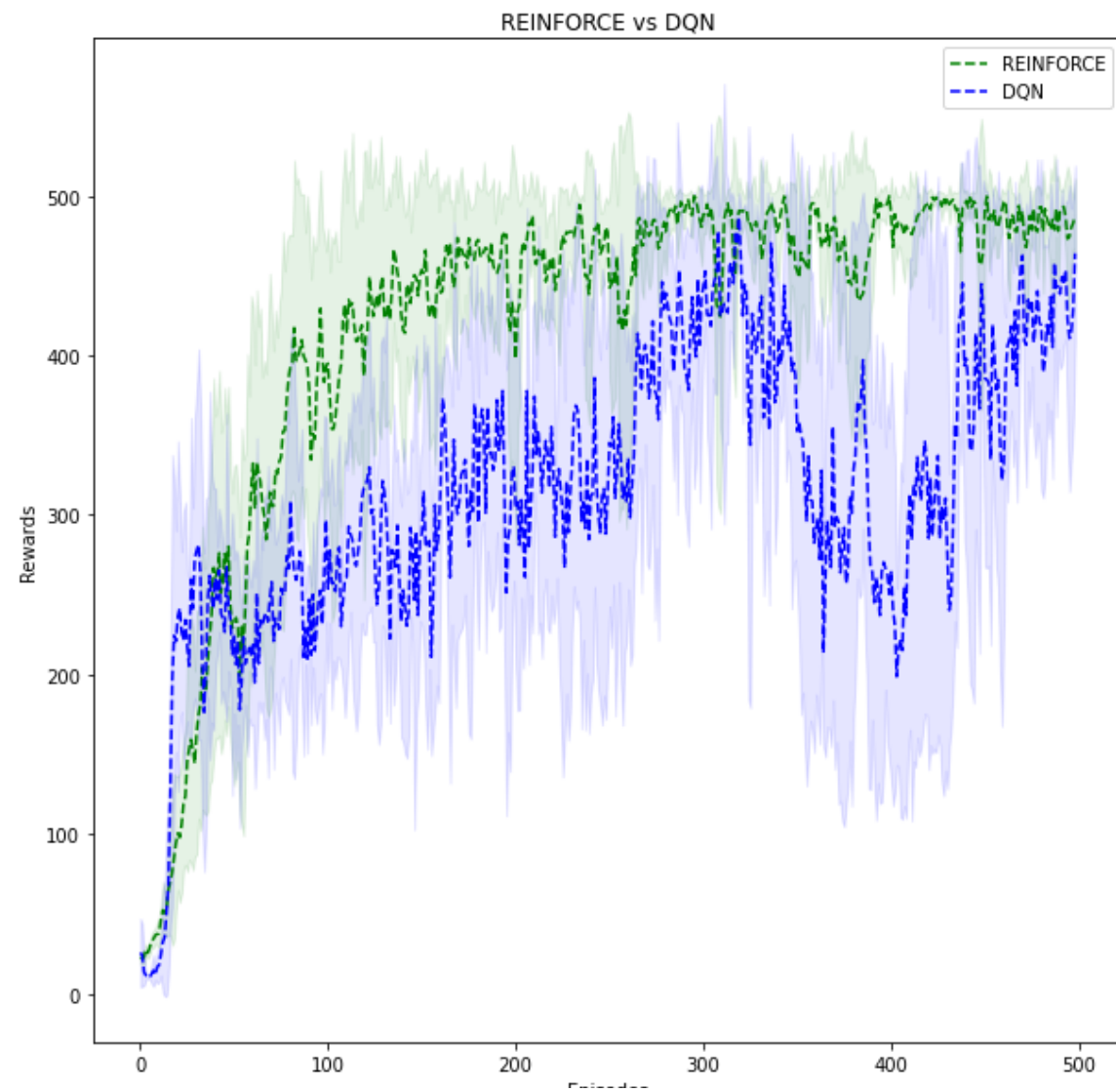
[Cart Position, Cart Velocity, Pole Angle, Pole Angular Velocity]

[Push cart to the left, Push cart to the right]

= to make pole remain upright







A2C

IDEA > Q-function 자리에 임의의 state function을 대입한다면?

$$\text{PROOF} > \forall \theta J_{\theta} \approx \int_{s_t, a_t} \nabla_{\theta} \ln P_{\theta}(a_t | s_t) \underbrace{Q(s_t, a_t)}_{\text{state function}} \underbrace{P_{\theta}(s_t, a_t)}_{\text{policy}} ds_t, da_t$$

$$\Rightarrow \int_{s_t, a_t} \nabla_{\theta} \ln P_{\theta}(a_t | s_t) \cdot A \cdot P_{\theta}(a_t | s_t) P(s_t) ds_t, da_t$$

$$\checkmark \nabla_{\theta} \ln P_{\theta}(a_t | s_t) = \nabla_{\theta} P_{\theta}(a_t | s_t) / P_{\theta}(a_t | s_t)$$

$$= \int_{s_t} \nabla_{\theta} \int_{a_t} P_{\theta}(a_t | s_t) da_t \cdot P(s_t) ds_t$$

$$\checkmark \nabla_{\theta} \cdot \perp = \text{ZERO}$$

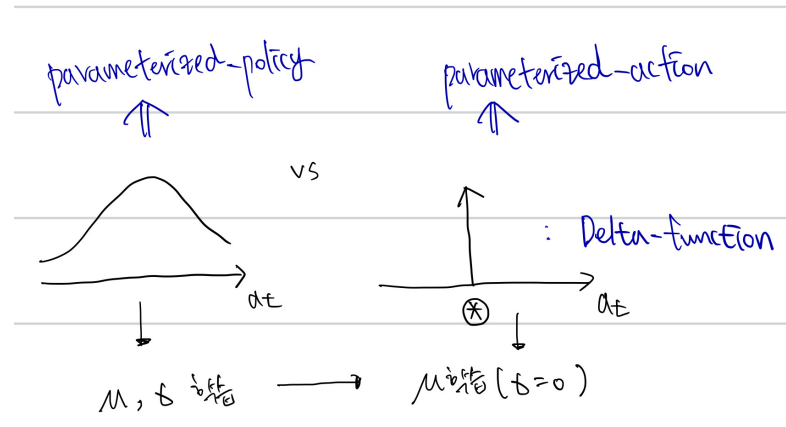
CONCLUSION > Q 대신 Q-V 함수 사용 : to lower variance

A2C

$$\begin{aligned}
 \nabla_{\theta} J_{\theta} &\approx \sum_{t=0}^{\infty} \int_{\mathbf{s}_t, \mathbf{a}_t} \nabla_{\theta} \ln p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) [Q(\mathbf{s}_t, \mathbf{a}_t) - V(\mathbf{s}_t)] p_{\theta}(\mathbf{s}_t, \mathbf{a}_t) d\mathbf{s}_t, \mathbf{a}_t \\
 &= \sum_{t=0}^{\infty} E [\nabla_{\theta} \ln p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) [\underbrace{Q(\mathbf{s}_t, \mathbf{a}_t) - V(\mathbf{s}_t)}_{\text{given}}]] \quad Q(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\mathbf{s}_{t+1}} [R_t + \gamma V(\mathbf{s}_{t+1}) | \mathbf{s}_t] \\
 &= \sum_{t=0}^{\infty} E_{\mathbf{s}_t, \mathbf{a}_t} [\nabla_{\theta} \ln p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \mathbb{E}_{\mathbf{s}_{t+1}} [R_t + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t) | \mathbf{s}_t]] \quad "V \rightarrow V_w" \\
 &\downarrow \\
 &= \sum_{t=0}^{\infty} \int_{\mathbf{s}_t, \mathbf{a}_t} \nabla_{\theta} \ln p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \int_{\mathbf{s}_{t+1}} (R_t + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) d\mathbf{s}_{t+1} p_{\theta}(\mathbf{s}_t, \mathbf{a}_t) d\mathbf{s}_t, \mathbf{a}_t \\
 &= \sum_{t=0}^{\infty} \int_{\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}} \nabla_{\theta} \ln p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) (\underbrace{R_t + \gamma V(\mathbf{s}_{t+1}) - V(\mathbf{s}_t)}_{\text{sample}}) p(\mathbf{s}_t) p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) d\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}
 \end{aligned}$$

(on-policy : update에 사용될 N개의 sample은 폐기해야함. → 무린 theta에서 뽑은 sample 이어서)

DDPG



$$J_{\theta} = E[G_t] = \int_{s_0:a_{\infty}} G_t \cdot P(s_0:a_{\infty}) ds_0:a_{\infty} = \int_{s_0:a_{\infty}} G_t \cdot P(a_0:a_{\infty} | s_0) da_0:a_{\infty} P(s_0) ds_0 = \int_{s_0} V(s_0) P(s_0) ds_0$$

$$V(s_0) = \int_{a_0:a_{\infty}} G_t \cdot P(a_0, \dots, a_{\infty} | s_0) da_0:a_{\infty} = \int_{a_0} \int_{s_1:a_{\infty}} G_t \cdot P(s_1:a_{\infty} | s_0, a_0) ds_1:a_{\infty} P(a_0 | s_0) da_0 = \int_{a_0} Q(s_0, a_0) P(a_0 | s_0) da_0$$

$$= Q(s_0, a_{\theta,0})$$

$$\text{cf } Q(s_0, a_0) = \int_{s_1:a_{\infty}} (R_0 + \gamma G_1) P(s_1:a_{\infty} | s_0, a_0) ds_1:a_{\infty} = \int_{s_1,a_1} \int_{s_2:a_{\infty}} (R_0 + \gamma G_1) P(s_2:a_{\infty} | s_1, a_1) P(a_1 | s_1) P(s_1 | s_0, a_0) ds_1, a_1 ds_2:a_{\infty}$$

$$= \int_{s_1,a_1} (R_0 + \gamma Q(s_1, a_1)) P(a_1 | s_1) P(s_1 | s_0, a_0) ds_1, a_1$$

$$Q(s_0, a_{\theta,0}) = R(s_0, a_{\theta,0}) + \int_{s_1} \gamma Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 \quad (\text{DDPG ver})$$

DDPG

$$\textcircled{1} \quad \nabla_{\theta} J_{\theta} = \int_{s_0} \nabla_{\theta} V(s) P(s_0) ds_0 = \int_{s_0} \nabla_{\theta} Q(s_0, a_{\theta,0}) P(s_0) ds_0$$

$$\nabla_{\theta} Q(s_0, a_{\theta,0}) = \nabla_{a_{\theta,0}} R(s_0, a_{\theta,0}) \cdot \nabla_{a_{\theta,0}} + \int_{s_1} \gamma Q(s_1, a_{\theta,1}) \nabla_{a_{\theta,0}} P(s_1 | s_0, a_{\theta,0}) \cdot \nabla_{a_{\theta,0}} ds_1 + \int_{s_1} \gamma \nabla_{\theta} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1$$

$$\nabla_{a_{\theta,0}} \cdot \nabla_{a_{\theta,0}} (R(s_0, a_{\theta,0}) + \int_{s_1} \gamma Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1) + \int_{s_1} \gamma \nabla_{\theta} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1$$

$$\nabla_{a_{\theta,0}} \cdot \nabla_{a_{\theta,0}} Q(s_0, a_{\theta,0}) + \int_{s_1} \gamma \nabla_{\theta} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1$$

$$\textcircled{2} \quad \nabla_{\theta} J_{\theta} = \int_{s_0} \nabla_{\theta} a_{\theta,0} \nabla_{a_{\theta,0}} Q(s_0, a_{\theta,0}) P(s_0) ds_0 + \int_{s_0} \int_{s_1} \nabla_{\theta} a_{\theta,1} \nabla_{a_{\theta,1}} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 P(s_0) ds_0$$

$$+ \int_{s_0} \int_{s_1} \int_{s_2} \nabla_{\theta} a_{\theta,2} \nabla_{a_{\theta,2}} Q(s_2, a_{\theta,2}) P(s_2 | s_1, a_{\theta,1}) ds_2 P(s_1 | s_0, a_{\theta,0}) ds_1 P(s_0) ds_0$$

$$= \sum_{t=0}^{\infty} \int_{s_0: s_t} \nabla_{\theta} a_{\theta,t} \nabla_{a_{\theta,t}} Q(s_t, a_{\theta,t}) P(s_0) P(s_1 | s_0, a_{\theta,0}) \cdot P(s_t | s_{t-1}, a_{\theta,t-1}) ds_0: s_t$$

Pendulum-v0

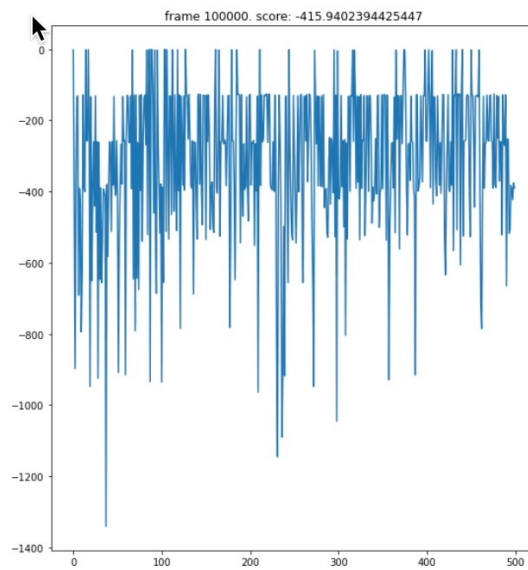
$[\cos(\theta), \sin(\theta), \dot{\theta}]$

$[\text{torque}]$

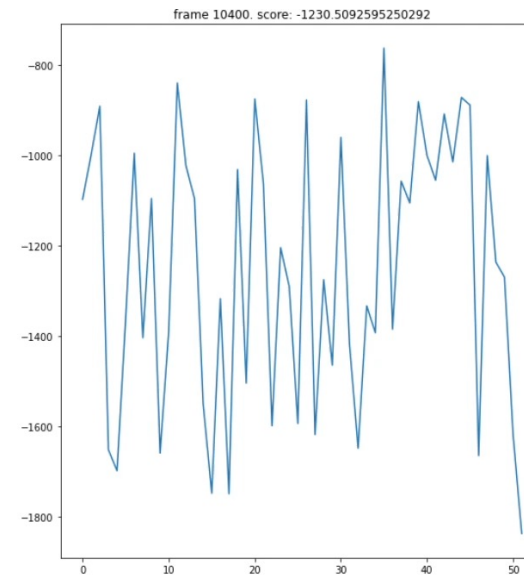
= *swing the pendulum up so it stays upright*



Results

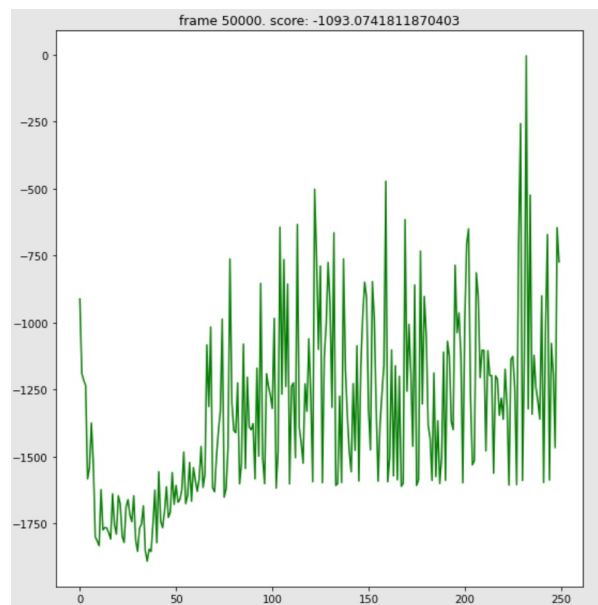


A2C

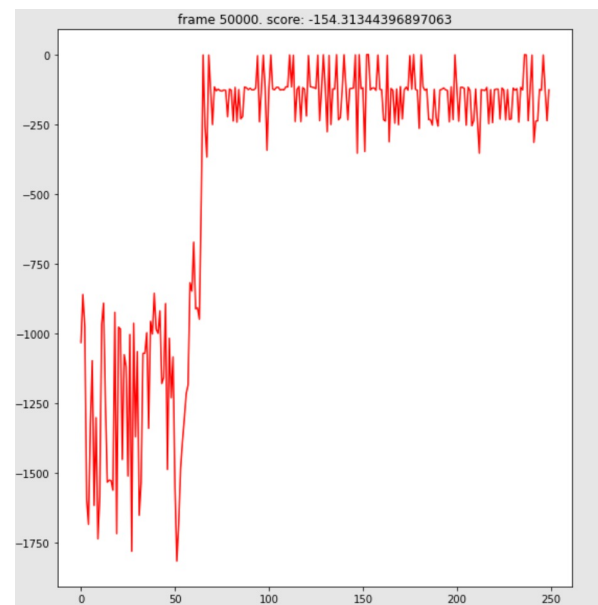


DDPG

Results



A2C

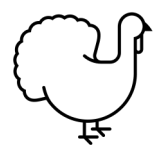


DDPG



Further Study

- 1) Off Policy Algorithm -> Different Policy (but almost similar)
- 2) Several Agents learning with totally different Algorithms / Policies
- 3) Sharing Results / Parameters of Learning -> Find out Optimal Policy / Agent
- 4) Apply to Sports Environment or Economic Model Simulator



감사합니다

