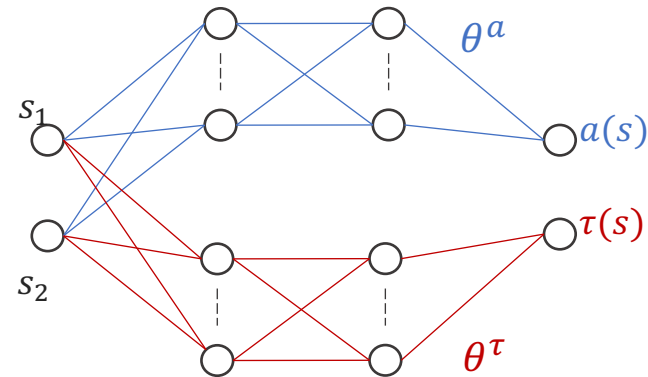


Weekly Report

M2 Ibuki Takeuchi

- Last week
 - Preparation for seminar
 - Calculate policy gradient



- Analytical calculation
 - Use dummy parameter θ to discuss what θ^a, θ^r have in common
 - θ^π is a tuple of θ^a, θ^r
 - $\nabla_\theta V^{\theta^\pi}(s) = \nabla_\theta [r(s, \pi(s|\theta^\pi)) + \gamma V^{\theta^\pi}(s'(\theta^\pi))]$
 $= \nabla_\theta r(s, \pi(s|\theta^\pi)) + \gamma \nabla_\theta \{V^{\theta^\pi}(s'(\theta^\pi))\}$
 $= \nabla_\theta r(s, \pi(s|\theta^\pi)) + \gamma \nabla_\theta s'(\theta^\pi) \nabla_{s'} V^{\theta^\pi}(s')|_{s'=s'(\theta^\pi)}$

$$+ \gamma \nabla_\theta V^{\theta^\pi}(s')|_{s'=s'(\theta^\pi)}$$

1st step's term

recursive term

- Detail calculation of $\nabla_{\theta^\tau} V^{\theta^\pi}(s)$
 - $s'(\theta^\pi) = s'(\theta^\pi|s)$: next state (τ second after)
 - $\nabla_{\theta^\tau} r(s, \pi(s|\theta^\pi)) = \nabla_{\theta^\tau} \tau(s|\theta^\tau) \{ \underbrace{-s'(\theta^\pi)^T Q s'(\theta^\pi) - a(s|\theta^\pi)^T R a(s|\theta^\pi) + \lambda}_{\nabla_\tau r(s, \pi(s|\theta^\pi))} \}$
 - $\nabla_{\theta^\tau} s'(\theta^\pi) = \nabla_{\theta^\tau} \tau(s|\theta^\tau) \{ \underbrace{f(s'(\theta^\pi)) + g(s'(\theta^\pi))a(s|\theta^a)}_{\nabla_\tau s'(\theta^\pi)} \}$
 - $\nabla_{\theta^\tau} V^{\theta^\pi}(s) = \nabla_{\theta^\tau} \tau(s|\theta^\tau) \{$
 $\quad - s'(\theta^\pi)^T Q s'(\theta^\pi) - a(s|\theta^\pi)^T R a(s|\theta^\pi) + \lambda$
 $\quad + \gamma (f(s'(\theta^\pi)) + g(s'(\theta^\pi))a(s|\theta^a)) \nabla_{s'} V^{\theta^\pi}(s')|_{s'=s'(\theta^\pi)} \}$
 $\quad + \gamma \nabla_{\theta} V^{\theta^\pi}(s')|_{s'=s'(\theta^\pi)}$
- Considering the equation
 - Is there region in parameter space where policy gradient steep?

- Strategy
 - If there is the reason, why the policy gradient suddenly changed, in $\nabla_{\theta^{\tau}} V^{\theta^{\pi}}(s), \nabla_{\theta^a} V^{\theta^{\pi}}(s)$
 - Ahead for master thesis(?)
 - If not, it means the reason above is in DDPG's approximation
 - For master thesis, compare $\nabla_{\theta^{\tau}} V^{\theta^{\pi}}(s), \nabla_{\theta^a} V^{\theta^{\pi}}(s)$ for theoretical hyper parameter settings