

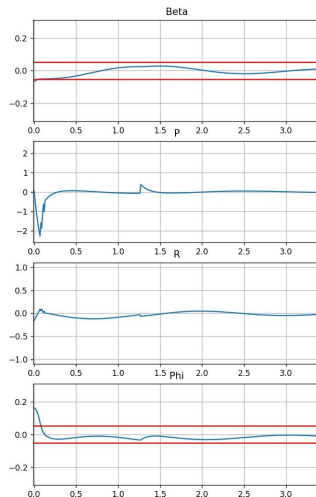
Meeting

Po Hsun Wu

September, 2022

Progress report

- Writing conference paper
- Study PPO and TRPO algorithm



TRPO algorithm

- TRPO optimal problem

$$\begin{aligned} \max_{\theta} \quad & \hat{\mathbb{E}}_t \left[\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t \right] \\ \text{subject to} \quad & \hat{\mathbb{E}}_t \left[\text{KL} \left[\pi_{\theta_{old}}(\cdot|s_t), \pi_{\theta}(\cdot|s_t) \right] \right] \leq \delta \end{aligned} \quad (1)$$

- Trust Region method is to solving nonlinear problem

$$\begin{aligned} \min \quad & m_k(p) = f_k + g_k^T p + \frac{1}{2} p^T B_k p \\ \text{subject to} \quad & \|p\| \leq \Delta_k \end{aligned} \quad (2)$$

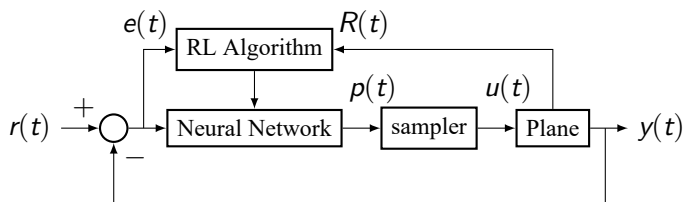
PPO algorithm

- PPO(Proximal Policy Optimization) is upgrade from TRPO(Trust Region Policy Optimization)
- PPO optimal problem

$$\min_{\theta} \hat{\mathbb{E}}_t \left[\min \left(\frac{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_{\theta_{\text{old}}}(\mathbf{a}_t | \mathbf{s}_t)} \hat{A}_t, \text{clip} \left(\frac{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_{\theta_{\text{old}}}(\mathbf{a}_t | \mathbf{s}_t)}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t \right) \right] \quad (3)$$

clip function is to set the policy change inside $[1 - \epsilon, 1 + \epsilon]$

Architecture



$r(t)$: reference
 $y(t)$: state
 $e(t)$: error
 $R(t)$: reward
 $p(t)$: probability
 of actions
 $u(t)$: control