

$${}^{hl}\mathcal{E}(\pi_l)\pi_l\pi_h$$

$$\eta(\pi_h)=V_h(s_0^h)=E_{s_0^h,a_0^h,\dots}\left[\sum_{t=0,k,2k,\dots}\gamma_h^{t/k}r_h(s_t^h,a_t^h)\right],where s_0^h\sim\rho_0^h(s_0^h),a_t^h\sim\pi_h(s_t^h),s_{t+1}\sim P(s_{t+1}^h|s_t^h,a_t^h,\mathcal{E}(\pi_l))$$

(1)

$$\begin{array}{l} \pi_h\pi_l\\ \frac{\eta(\pi_h)RR}{\frac{\gamma_hV_h(s_{t+k}^h)-V_h(s_t^h)}{k}\pi_l}\\ \text{?}V_\pi(s_0)\text{?}JV_\pi(s_0)RE_{s,a\sim\pi}[R(s,a)]\\ R\\ \nabla J(\theta)=\\ \nabla E_{s\sim\pi}[\sum_a\pi(s,a)R(s,a)]=\\ \nabla E_{s,a\sim\pi}[\hat{R}(s,a)].\\ l(s_{t+i}^l,a_{t+i}^l)|_{i=0,1,\dots,k-1}=\\ \frac{V_h(s_{t+k}^h)-V_h(s_t^h)}{k}.\\ E_{s,a\sim\pi}[R(s,a)]E_{s,a\sim\pi_l,\pi_h}[R_l(s_l,a_l)]\pi_h\pi_h\\ \text{low}E_{s^l,a^l\sim\pi_l,\pi_h}[R_l(s^l,a^l)]=\\ \frac{1}{k}E_{s^h,a^h\sim\pi_l,\pi_h}[\gamma_hV_h(s_{t+k}^h)-\\ V_h(s_t^h)].\\ A(s,a)\\ {}_t,A_t)=\\ Q(S_t,A_t)-\\ V(S_t)=\\ R(S_t,A_t)+\\ \gamma V(S_{t+1})-\\ V(s).\\ {}_t,A_t)=\\ 0,\forall t\neq\\ t_{end}.\\ {}_t,A_t)=\\ \gamma V(S_{t+1})-\\ V(s).\\ expectationeq:sparseadvantageTRPOE_{s^l,a^l\sim\pi_l,\pi_h}[R_l(s^l,a^l)]=\\ \frac{1}{k}E_{s^h,a^h\sim\pi_l,\pi_h}[A^h(s_t^h,a_t^h)].\\ \tilde{\eta}(\pi_h)\pi_l\tilde{\pi}_l\pi_h\text{?}\mathbf{Lemma}\\ \mathbf{1}\\ \tilde{\eta}(\pi_h)=\\ \eta(\pi_h)+\\ E_{s_0^h,a_0^h,\dots\sim\pi_h,\mathcal{E}(\pi_l)}\left[\sum_{t=0,k,2k,\dots}\gamma_h^{t/k}A_{\pi_h}(s_t^h,a_t^h)\right].\\ environmentadvantage\\ \rho\\ \rho_{\pi_h}(s^h)=\\ \sum_{t=0,k,2k,\dots}\gamma_h^{t/k}P(s_t^h=\\ s|\mathcal{E}(\pi_l))\\ \rho_{\pi_h}\rho_{\pi_h}\\ environmentadvantageTRPO22\sum_{s^h}\tilde{\rho}_{\pi_h}(s^h)\sum_{a^h}\pi_h(a^h|s^h)A_{\pi_h}(s^h,a^h).\\ {}_h=\\ {}_0^h)=\\ P(s_1^h=\\ s^h)=\\ \ddots\\ \ddot{P}(s_i^h=\\ s^h)=\\ P(s_i^h=\\ s^h)=\\ objective\tilde{\rho}_{\pi_h}(s^h)=\\ \frac{1}{1-\gamma_h}P(s_i^h=\\ s^h)\\ environmentadvantage\tilde{\eta}(\pi_h)=\\ \eta(\pi_h)+\\ \frac{1}{1-\gamma_h}E_{s^h,a^h\sim\pi_h,\mathcal{E}(\pi_l)}\left[A_{\pi_h}(s^h,a^h)\right].\\ expectationisadvantageeq:accurateobjectiveexpectationformTRPO\eta(\pi_h)\\ \pi_h\pi_l\pi_l\pi_h\eta(\pi_h)\\ \pi_h\pi_l\eta(\pi_h)\pi_l\pi_h\\ \mathbf{1}\\ \tilde{\eta}(\pi_h)=\\ \eta(\pi_h)+\end{array}$$