Reparameterized Policy Learning for Multimodal Trajectory Optimization

Zhiao Huang¹ Litian Liang¹ Zhan Ling¹ Xuanlin Li¹ Chuang Gan²³ Hao Su¹

¹UC San Diego ²MIT-IBM Watson AI Lab ³UMass Amherst. Correspondence to: Zhiao Huang <z2huang@ucsd.edu>.

ICML2023

ParN₂C 模式识别与神经计算研究组 PAttern Recognition and NEural Computing

Variational Inference



$$\mathbb{E}_{q(z)}[\log q(z) - \log p(z, x) + \log p(x)] \ge 0 \Rightarrow \log \underbrace{p(x)}_{\text{evidence}} \ge \underbrace{\mathbb{E}_{q(z)}[\log p(z, x) - \log q(z)]}_{\text{ELBO}}$$

RL as probabilistic inference



$$p(\tau|O_{1:T}) = p(\tau, O_{1:T}) p(O_{1:T}) = p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) p(O_t|s_t, a_t) \prod_{t=1}^{T} p(O_t)$$

 $p(O_t)$ are the optimal variable prior which can be considered as a constant

Optimal variable
$$O_t: p(O_t|s_t, a_t) = \exp(r(s_t, a_t))$$

 $p(\tau|O_{1:T}) \propto p(\tau, O_{1:T}) = \left[p(s_1) \prod_{t=1}^T p(a_t) p(s_{t+1}|s_t, a_t) \right] \exp\left(\sum_{t=1}^T r(s, a)\right)$

$$\hat{p}^{\pi}(\tau) = p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) \frac{\pi(a_t|s_t)}{p(a_t)} = \left[p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) \right] \prod_{t=1}^{T} \frac{\pi(a_t|s_t)}{p(a_t)}$$

$$p(\tau|O_{1:T}) \propto p(\tau, O_{1:T}) = \left[p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) \right] \exp\left(\sum_{t=1}^{T} r(s, a)\right)$$

$$egin{aligned} \min \mathrm{D}_{\mathit{KL}}ig(\hat{p}^{\,\pi}(au)||p(au|O_{1:T})ig) &\Leftrightarrow \max - \mathbb{E}_{\, au \sim \, \hat{p}(au, \pi)}ig[\log \hat{p}\left(au, \pi
ight) - \log p\left(au, O_{1:T}
ight)ig] \ &= \max \, \mathbb{E}_{\, au \sim \, \hat{p}(au, \pi)}igg[\sum_{t=1}^Tig(r(s_t, a_t) - \log \pi\left(a_t|s_t
ight) + \log p\left(a_t
ight) ig)ig] \ &\Leftrightarrow \max \, \sum_{t=1}^T \, \mathbb{E}_{s_t \sim \, \hat{p}(s_t), a_t \sim \, \pi\left(\cdot \mid s_t
ight)}ig[r(s_t, a_t) - \log \pi\left(a_t|s_t
ight)ig] \end{aligned}$$

or $\log p(O_{1:T}) \ge \mathbb{E}_{\tau \sim \pi} \left[\log p(\tau, O_{1:T}) - \log \hat{p}^{\pi}(\tau) \right] = \mathbb{E}_{\tau \sim \pi} \left[\log p(O_{1:T}|\tau) + \log p(\tau) - \log \hat{p}^{\pi}(\tau) \right]$

Motivation

ParNeC 模式识别与神经计算研究组 PAttern Recognition and NEural Computing



As the distribution of rewards (or Q function) can be multimodal, it is hard to evade local optima if the policy is modeled as a unimodal distribution

Method



Reparameterized Policy $\pi_{\theta}(z,\tau) = p(s_1)\pi_{\theta}(z|s_1) \prod_{t=1}^{T} p(s_{t+1}|s_t,a_t)\pi_{\theta}(a_t|z,s_t)$

 $\log p(O_{1:T}) \ge \mathbb{E}_{q(z)} \left[\log p(\tau, O_{1:T}) - \log \hat{p}^{\pi}(\tau) \right] = \mathbb{E}_{q(z)} \left[\log p(O_{1:T} | \tau) + \log p(\tau) - \log \hat{p}^{\pi}(\tau) \right]$

$$\begin{split} &\log p(O) \\ &= \underbrace{E_{z,\tau \sim \pi_{\theta}} \left[\log p_{\phi}(O, z, \tau) - \log \pi_{\theta}(z, \tau) \right]}_{\text{ELBO}} \\ &+ D_{KL}(\pi_{\theta}(z, \tau)) || p_{\phi}(z, \tau | O)) \\ &\geq E_{z,\tau \sim \pi_{\theta}} \left[\log p_{\phi}(O, \tau, z) - \log \pi_{\theta}(z, \tau) \right] \quad p_{\phi}(O, z, \tau) = p(O|z, \tau) p(z|\tau) p(\tau) = p(O|\tau) p(z|\tau) p(\tau) \\ &= E_{z,\tau \sim \pi_{\theta}} \left[\log p(O, \tau) + \log p_{\phi}(z|\tau) - \log \pi_{\theta}(z, \tau) \right] \\ &= E_{z,\tau} \left[\underbrace{\log p(O|\tau)}_{\text{reward}} + \underbrace{\log p(\tau)}_{\text{prior}} + \underbrace{\log p_{\phi}(z|\tau)}_{\text{cross entropy}} - \underbrace{\log \pi_{\theta}(z, \tau)}_{\text{entropy}} \right] \quad \text{auxiliary distribution} \end{split}$$

Method



$$E_{z,\tau} \left[\underbrace{\log p(O|\tau)}_{\text{reward}} + \underbrace{\log p(\tau)}_{\text{prior}} + \underbrace{\log p_{\phi}(z|\tau)}_{\text{cross entropy}} - \underbrace{\log \pi_{\theta}(z,\tau)}_{\text{entropy}} \right]$$

reward term:
$$p(O|s,a) = \exp(r(s,a)) / T$$

prior term: a constant

cross entropy term: model the posterior of z for τ sampled from π and distinguish different trajectories by z $I(\tau,z) = H(z) - H(z|\tau) = \mathbb{E}_{\tau,z} [\log p(z|\tau) - \log p(\tau)] = \mathbb{E}_{\tau,z} \Big[\log \frac{p(z|\tau)}{p_{\phi}(z|\tau)} + \log p_{\phi}(z|\tau) - \log p(\tau) \Big] \ge \mathbb{E}_{\tau,z} [\log p_{\phi}(z|\tau) - \log p(\tau)]$

entropy term: encourage exploration like maximum entropy RL

$$\log p_{\phi}(z|\tau) = \sum_{t>0} \log p(z|s_t, a_t)$$
final reward:
$$\underbrace{\frac{R(s_t, a_t)}{\mathcal{T}}}_{r_t} \underbrace{-\alpha \log \pi_{\theta}(a_t|s_t, z) + \beta \log p_{\phi}(z|s_t, a_t)}_{r'_t}$$

Method



Model-based RL

 $s_{t} = f_{\psi}(o_{t})$ $r_{t} = R_{\psi}(s_{t}, a_{t})$ $Q_{t} = Q_{\psi}(s_{t}, a_{t})$ $s_{t+1} = h_{\psi}(s_{t}, a_{t})$ (A) Reparameterized Policy $V_{est}(o_{t_{0}}, z) \approx \gamma^{K}(Q_{t_{0}+K} + r'_{t_{0}+K}) + \sum_{t=t_{0}}^{t_{0}+K-1} \gamma^{t-t_{0}}(r_{t} + r'_{t})$ $L_{\psi}(\tau) = \sum_{t=t_{0}}^{t_{0}+K-1} L_{1} ||s_{t+1} - \mathbf{ng}(f_{\psi}(o_{t+1}))||^{2} + L_{2}(r_{t} - r_{t}^{gt})^{2}$ $+ L_{3}(Q_{t} - \mathbf{ng}(r_{t}^{gt} + \gamma V_{est}(o_{t+1}, z)))^{2} \qquad (4)$ (C) Exploration Bonus $-\log p(s)$



Experiment



Algorithm 1 Model-based Reparameterized Policy Gradient

Input: $p_{\phi}, \pi_{\theta}, h_{\psi}, R_{\psi}, f_{\psi}, Q_{\psi}$ and an optional density estimator g_{θ}

Initialize p_{ϕ}, π_{θ} , construct the replay buffer \mathcal{B} .

while time remains do

Sample start state o_1 and encode it as $s_1 = f_{\psi}(o_1)$. Select z from $\pi_{\theta}(z|s_1)$.

Execute the policy $\pi_{\theta}(a|s, z)$ and store transitions into the replay buffer \mathcal{B} .

Sample a batch of trajectory segment of length $K \{\tau_{t:t+K}^i, z\}$ from the buffer \mathcal{B} .

Optional: update and estimate the density estimator g_{θ} and relabel transitions with the negative density as the intrinsic reward.

Optimize ψ using Equation 4.

Optimize $\pi_{\theta}(a|s, z)$ with gradient descent to maximize the value estimate in Equation 3 for s, z sampled from the buffer.

Optimize $\pi_{\theta}(z|s_1)$ with policy gradient to maximize $V_{\text{estimate}}(s_1, z) - \alpha \log \pi_{\theta}(z|s_1)$ for s_1 sampled from the buffer. Optimize α, β if necessary.

end while





Can multimodal policies help escape local optima?

Figure 4. Illustrative experiment on continuous bandit

Can multimodal policies accelerate exploration?



Figure 5. Illustrative experiment on 2D maze navigation problem

Experiment



Figure 6. Results on dense reward tasks with local optima (exploration disabled)

Figure 7. Results on sparse reward tasks

PBRL

$$P_{\psi}[\sigma^{1} \succ \sigma^{0}] = \frac{\exp \sum_{t} \hat{r}_{\psi}(\mathbf{s}_{t}^{1}, \mathbf{a}_{t}^{1})}{\sum_{i \in \{0,1\}} \exp \sum_{t} \hat{r}_{\psi}(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i})}$$
$$\mathcal{L}^{\text{Reward}} = -\underset{(\sigma^{0}, \sigma^{1}, y) \sim \mathcal{D}}{\mathbb{E}} \left[y(0) \log P_{\psi}[\sigma^{0} \succ \sigma^{1}] + y(1) \log P_{\psi}[\sigma^{1} \succ \sigma^{0}] \right]$$



$$egin{aligned} \log p\left(O_{1:T}
ight) &\geq \mathbb{E}_{ au \sim \pi} ig[\log p\left(au, O_{1:T}
ight) - \log \hat{p}^{\pi}\left(au
ight)ig] \ &= \mathbb{E}_{ au \sim \pi} ig[\log p\left(O_{1:T}| au
ight) + \log p\left(au
ight) - \log \hat{p}^{\pi}\left(au
ight)ig] \ &= \sum_{t=1}^T \mathbb{E}_{s_t \sim \hat{p}^{\pi}(s_t), a \sim \pi\left(\cdot \mid s_t
ight)} ig[r(s_t, a_t) - \log \pi\left(a_t|s_t
ight)ig] \end{aligned}$$

最小化与较好轨迹的KL散度,最大化与较差轨迹的KL散度 $D_{KL}(\hat{p}^{\pi}(\tau)||p(\tau|O_{1:T})) \Leftrightarrow \max \sum_{t=1}^{T} \mathbb{E}_{s_t \sim \hat{p}^{\pi}(s_t), a_t \sim \pi(\cdot|s_t)} [r(s_t, a_t) - \log \pi(a_t|s_t)]$ $\Leftrightarrow \max \sum_{t=1}^{T} \mathbb{E}_{s_t \sim \hat{p}^{\pi}(s_t), a_t \sim \pi(\cdot|s_t)} [r(s_t, a_t)] + \sum_{t=1}^{T} \mathbb{E}_{s_t \sim \hat{p}^{\pi}(s_t)} [H(\pi(\cdot|s_t))]$

注意到
$$\hat{p}^{\pi}(\tau) = p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) \frac{\pi(a_t|s_t)}{p(a_t)} = \left[p(s_1) \prod_{t=1}^{T} p(a_t) p(s_{t+1}|s_t, a_t) \right] \prod_{t=1}^{T} \frac{\pi(a_t|s_t)}{p(a_t)},$$
其中第一项
与轨迹 τ 的转移动态一致, 故假设 $\hat{p}^{\pi}(s_t) = p^{\tau}(s_t)$

$$\begin{split} \mathrm{D}_{KL}\big(\hat{p}^{\,\pi}(\tau)||p(\tau_w|O_{1:T})\big) - \mathrm{D}_{KL}\big(\hat{p}^{\,\pi}(\tau)||p(\tau_l|O_{1:T})\big) &\Leftrightarrow \max \ \sum_{t=1}^T \mathbb{E}_{(s_t,a_t) \sim \tau_w} \frac{\pi(a_t|s_t)}{\pi_{\tau_w}(a_t|s_t)} [r(s_t,a_t)] + \sum_{t=1}^T \mathbb{E}_{(s_t,a_t) \sim \tau_w} \big[H\big(\pi(\cdot|s_t)\big)\big] \\ &- \sum_{t=1}^T \mathbb{E}_{(s_t,a_t) \sim \tau_t} \frac{\pi(a_t|s_t)}{\pi_{\tau_t}(a_t|s_t)} [r(s_t,a_t)] - \sum_{t=1}^T \mathbb{E}_{(s_t,a_t) \sim \tau_t} \big[H\big(\pi(\cdot|s_t)\big)\big] \end{split}$$

対于PBRL, 目标为最大化exp
$$\left(\sum_{t=1}^{T_{w}} \mathbb{E}_{\tau \sim p(\tau_{w})}[r(s_{t},a_{t})] - \sum_{t=1}^{T_{1}} \mathbb{E}_{\tau \sim p(\tau_{w})}[r(s_{t},a_{t})]\right)$$

对比而言, 应该最大化exp $\left(-D(\hat{p}^{\pi}(\tau)||p(\tau_{w}|O_{1:T})) + D(\hat{p}^{\pi}(\tau)||p(\tau_{1}|O_{1:T}))\right)$
$$= \exp\left(\sum_{t=1}^{T_{w}} \mathbb{E}_{\tau \sim p(\tau_{w})} \frac{\pi(a_{t}|s_{t})}{\pi_{\tau_{w}}(a_{t}|s_{t})}[r(s_{t},a_{t})] - \sum_{t=1}^{T_{1}} \mathbb{E}_{\tau \sim p(\tau_{1})} \frac{\pi(a_{t}|s_{t})}{\pi_{\tau_{t}}(a_{t}|s_{t})}[r(s_{t},a_{t})]\right)$$

由于PBRL中较差轨迹和较好轨迹都是以前策略采样动作交互得到的,那么可以在策略采样动作记录 $\log \pi_b(a|s)$, 在学习奖励模型的时候即可计算重要性采样比 $\frac{\pi(a|s)}{\pi_b(a|s)} = \exp(\log \pi(a|s) - \log \pi_b(a|s))$,作为不同状态动作对下的奖励的权重





$$L(\phi) = -\log\left(\frac{\exp\left(\sum r_{\phi}(s_{t}^{1}, a_{t}^{1})\right)}{\exp\left(\sum r_{\phi}(s_{t}^{1}, a_{t}^{1})\right) + \exp\left(\sum r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}\right) \quad \frac{\partial L}{\partial r_{\phi}(s_{t}^{1}, a_{t}^{1})} = -\left(1 - \frac{\exp\left(\sum r_{\phi}(s_{t}^{1}, a_{t}^{1})\right) + \exp\left(\sum r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}{\exp\left(\sum r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}\right) \quad \frac{\partial L}{\partial r_{\phi}(s_{t}^{2}, a_{t}^{2})} = 1 - \frac{\exp\left(\sum r_{\phi}(s_{t}^{1}, a_{t}^{1})\right) + \exp\left(\sum r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}{\exp\left(\sum r_{\phi}(s_{t}^{1}, a_{t}^{1})\right) + \exp\left(\sum r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}\right)$$

$$L(\phi) = -\log\left(\frac{\exp\left(\sum A(s_{t}^{1}, a_{t}^{1})r_{\phi}(s_{t}^{1}, a_{t}^{1})\right)}{\exp\left(\sum A(s_{t}^{1}, a_{t}^{1})r_{\phi}(s_{t}^{1}, a_{t}^{1})\right) + \exp\left(\sum A(s_{t}^{2}, a_{t}^{2})r_{\phi}(s_{t}^{2}, a_{t}^{2})\right)}\right)$$

$$\frac{\partial L}{\partial r_{\phi}(s_{t}^{1},a_{t}^{1})} = -A(s_{t}^{1},a_{t}^{1}) \left(1 - \frac{\exp(\sum A(s_{t}^{1},a_{t}^{1})r_{\phi}(s_{t}^{1},a_{t}^{1}))}{\exp(\sum A(s_{t}^{1},a_{t}^{1})r_{\phi}(s_{t}^{1},a_{t}^{1})) + \exp(\sum A(s_{t}^{2},a_{t}^{2})r_{\phi}(s_{t}^{2},a_{t}^{2}))}\right)$$

$$rac{\partial L}{\partial r_{\phi}(s_{t}^{2},a_{t}^{2})} = A(s_{t}^{2},a_{t}^{2}) \left(1 - rac{\exp(\sum A(s_{t}^{1},a_{t}^{1})r_{\phi}(s_{t}^{1},a_{t}^{1}))}{\exp(\sum A(s_{t}^{1},a_{t}^{1})r_{\phi}(s_{t}^{1},a_{t}^{1})) + \exp(\sum A(s_{t}^{2},a_{t}^{2})r_{\phi}(s_{t}^{2},a_{t}^{2}))}
ight)$$

梯度 $\propto |A(s,a)(1-P)|$,最大化梯度以最大化每次更新能够获得的信息