

$$1) \sum_t [H(\theta | \epsilon_t, a_t) - H(\theta | s_{t+1}, \epsilon_t, a_t)]$$

$$I(x; y) = H(x) - H(x|y)$$

$$2) \therefore I(\theta; s_{t+1} | \epsilon_t, a_t)$$

$$= \mathbb{E}_{s_{t+1} \sim P(\cdot | \epsilon_t, a_t)} [D_{KL}[P(\theta | \epsilon_t, a_t, s_{t+1}) || P(\theta | \epsilon_t, a_t)]]$$

$$I(x; y) = \mathbb{E}_y [D_{KL}[P_{x|y} || P_x]]$$

$P(\theta | \epsilon_t)$

info gain

$$3) r'(s_t, a_t, s_{t+1}) = r(s_t, a_t) + \eta D_{KL}[P(\theta | \epsilon_t, a_t, s_{t+1}) || P(\theta | \epsilon_t)]$$

$$4) P(\theta | \epsilon_t, a_t, s_{t+1}) = \frac{P(\theta, \epsilon_t, a_t, s_{t+1})}{P(\epsilon_t, a_t, s_{t+1})}$$

chain rule

$$= \frac{P(\theta) P(\epsilon_t | a_t) P(a_t | \epsilon_t) P(s_{t+1} | \epsilon_t, a_t, \theta)}{P(\epsilon_t) P(a_t | \epsilon_t) P(s_{t+1} | \epsilon_t, a_t)}$$

$$= \frac{P(\theta | \epsilon_t) P(s_{t+1} | \epsilon_t, a_t, \theta)}{P(s_{t+1} | \epsilon_t, a_t)}$$

$$P(\theta | \epsilon_t, a_t)$$

$$= P(\theta | \epsilon_t)$$

$$5) P(s_{t+1} | \epsilon_t, a_t) = \int_{\theta} P(s_{t+1} | \epsilon_t, a_t, \theta) P(\theta | \epsilon_t) d\theta$$

hard to compute

Calculating $P(\theta | D)$ is hard for given dataset D

\therefore Approximate using $q(\theta; \phi)$ by minimizing

$$D_{KL}[q(\theta; \phi) || P(\theta | D)]$$

6) Variational Lower bound derivation

$$KL[q(z) || P(z|x)] = \int q(z) \log \frac{q(z)}{P(z|x)}$$

$$\begin{aligned}
&= - \int q(z) \log \frac{p(z|x)}{q(z)} \\
&= - \left[\int q(z) \log \frac{p(x,z)}{q(z)} - \int q(z) \log p(x) \right] \\
&= - \int q(z) \log \frac{p(x,z)}{q(z)} + \log p(x) \int q(z) \\
&= -L + \log p(x).
\end{aligned}$$

$$\therefore L = \log p(x) - KL[q(z) || p(z|x)]$$

$$L \leq \log p(x) \quad \text{as } KL \geq 0$$

Lower bound L hits the log probability iff the approx distro is perfectly close to the posterior distro. Hence maximise L

$$L[q(\theta; \phi), D] = \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log p(D|\theta)] - D_{KL}[q(\theta; \phi) || p(\theta)]$$

7.) Thus instead of calculating info gain in eq 3) we compute approx to it

$$\therefore r'(s_t, a_t, s_{t+1}) = r(s_t, a_t) + \eta D_{KL}[q(\theta; \phi_{t+1}) || q(\theta; \phi_t)]$$

ϕ_{t+1} and ϕ_t are new and old params representing the agent's belief

8.) ϕ_t is parametrized using BNMs

$$p(y|x) = \int_{\theta} p(y|x; \theta) q(\theta; \phi) d\theta$$

$$9) \log p(D) = \int_{\theta} \underbrace{q(\theta; \phi) \log \frac{p(\theta, D)}{q(\theta; \phi)}}_{L[q(\theta; \phi), D]} d\theta + D_{KL}[q(\theta; \phi) || p(\theta; D)]$$

Implementation

$$11) q(\theta; \phi) = \prod_{i=1}^{|\theta|} \mathcal{N}(\theta_i | \mu_i; \sigma_i^2)$$

$$\phi = \{\mu, \sigma\} \quad \sigma \rightarrow \text{covariance diagonal matrix}$$

$$\sigma = \log(1 + e^J) \quad J \in \mathbb{R} \quad \text{as } \sigma > 0$$

Egn 6

$$L[q(\theta; \phi, D)] = \underbrace{\mathbb{E}_{\theta \sim q(\cdot, \phi)} [\log p(D|\theta)]}_{\text{(VLB) variational lower bound}} - \underbrace{D_{KL}[q(\theta; \phi) \| p(\theta)]}_{\text{KL divergence}}$$

Approximated through sampling

$$\mathbb{E}_{\theta \sim q(\cdot, \phi)} [\log p(D|\theta)] \approx \frac{1}{N} \sum_{i=1}^N \log p(D|\theta_i)$$

with N samples drawn according to $\theta \sim q(\cdot; \phi)$

Refer stochastic gradient variational Bayes (SGVB)

The optimization of VLB is done at regular intervals by sampling D . This is done to break up intratrajectory sample correlation

12) ϕ for eqn 7) is calculated as

$$\phi' = \underset{\phi}{\text{argmin}} \left[\underbrace{D_{KL}[q(\theta; \phi) \| q(\theta; \phi_{t-1})]}_{L(q(\theta; \phi), st)} - \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log (s_t | \epsilon_t, q_t; \theta)] \right]$$

Here, expectation over θ is replaced with samples $\theta \sim q(\cdot; \phi)$

13) To optimize 12) efficiently,

$$\text{compute } D_{KL}[q(\theta; \phi + \lambda \Delta \phi) \| q(\theta; \phi)]$$

$$\Delta \phi = H^{-1}(L) \nabla_{\phi} L(q(\theta; \phi), st)$$

$H(L)$ is Hessian of $L(q(\theta; \phi), st)$

Since q is fully factorized Gaussian.

$$14) D_{KL}[q(\theta; \phi) \| q(\theta; \phi')]$$

$$= \frac{1}{2} \sum_{i=1}^{|\theta|} \left[\left(\frac{\sigma_i}{\sigma'_i} \right)^2 + 2 \log \sigma'_i - 2 \log \sigma_i + \frac{(\mu_i - \mu'_i)^2}{\sigma_i'^2} \right] - \frac{|\theta|}{2}$$

KL divergence is quadratic (approx) and the log-likelihood term can be seen as locally linear compared to the curved KL term, we approximate H by only calculating it for the KL term $L_{KL}(q(\theta; \phi))$

$$15.) \frac{\partial^2 L_{KL}}{\partial \mu_i^2} = \frac{1}{\log^2(1+e^{\mu_i})} \quad \text{and} \quad \frac{\partial^2 L_{KL}}{\partial \beta_i^2} = \frac{2e^{2\beta_i}}{(1+e^{\beta_i})^2} \frac{1}{\log^2(1+e^{\beta_i})}$$

$$16.) D_{KL}[q(\theta; \phi + \lambda \Delta \phi) \| q(\theta; \phi)] \approx \frac{1}{2} \lambda^2 \nabla_{\phi}^T H^{-1}(L_{KL}) \nabla_{\phi}$$

Here $H^{-1}(L_{KL})$ is diagonal

Algo

For each epoch n do

for each timestep t in trajectory do

Generate action $a_t \sim \pi_{\alpha}(s_t)$

and sample $s_{t+1} \sim P(\cdot | E_t, a_t)$, get $r_t(s_t, a_t)$

Add (s_t, a_t, s_{t+1}) in FIFO buffer R

Compute $D_{KL}[q(\theta; \phi_{n+1}) \| q(\theta; \phi_{n+1})]$ by approximation

$\nabla^T H^{-1} \nabla$ from (16) for diagonal BNNs

Divide $D_{KL}[q(\theta; \phi_{n+1}) \| q(\theta; \phi_{n+1})]$ by median of previous KL divergences

Construct $r'(s_t, a_t, s_{t+1}) \leftarrow r(s_t, a_t)$

Minimize $D_{KL}[q(\theta; \phi_n) \| p(\theta)] + \eta D_{KL}[q(\theta; \phi_{n+1}) \| q(\theta; \phi_{n+1})]$
 $- E_{\theta \sim q(\cdot; \phi_n)} [\log S(D|\theta)]$ (6) with

D sampled randomly from R , leading to updated posterior $q(\theta; \phi_{n+1})$

Use rewards $\{r'(s_t, a_t, s_{t+1})\}$ to update π_{α} using any standard RL algo.

