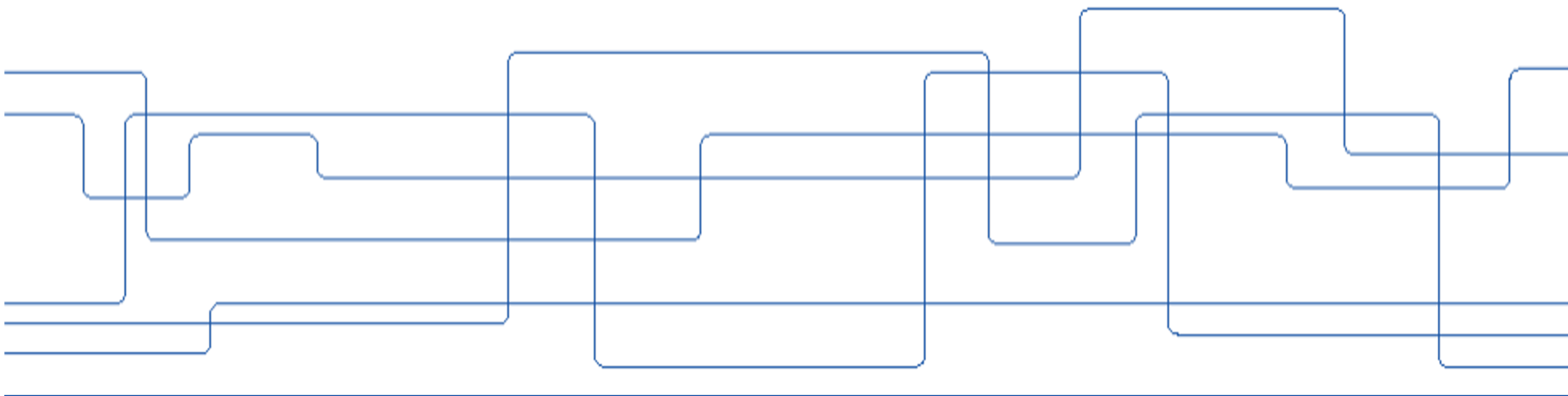




FDD3359 Reinforcement Learning Course

Offline RL

Ali Ghadirzadeh

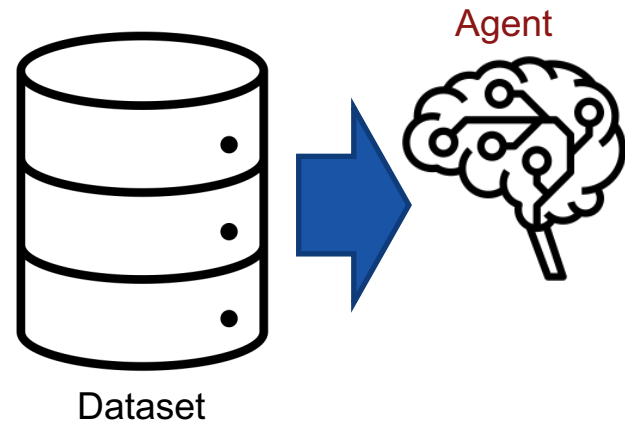
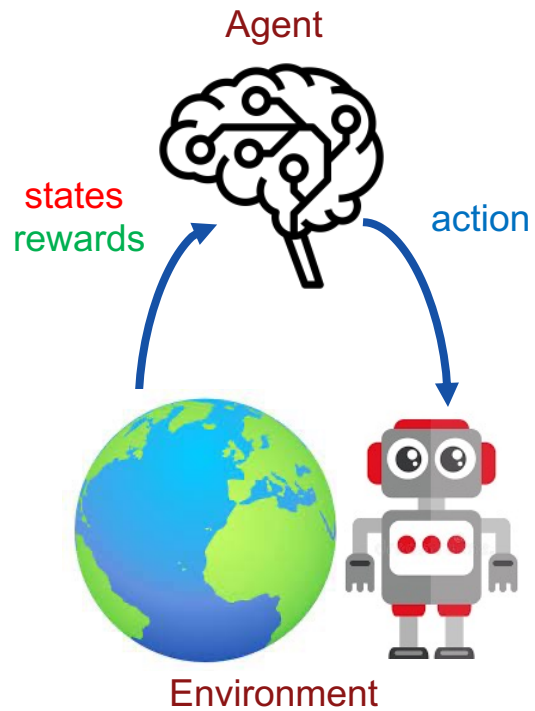


Offline RL?

A basic reinforcement learning agent \mathcal{A} interacts with its environment in discrete time steps. At each time t , the agent receives the current state s_t and reward r_t . It then chooses an action a_t from the set of available actions, which is subsequently sent to the environment. The environment moves to a new state s_{t+1} and the reward r_{t+1} associated with the *transition* (s_t, a_t, s_{t+1}) is determined. The goal of a reinforcement learning agent is to learn a *policy*: $\pi : A \times S \rightarrow [0, 1]$, $\pi(a, s) = \Pr(a_t = a \mid s_t = s)$ which maximizes the expected cumulative reward.



Offline RL?

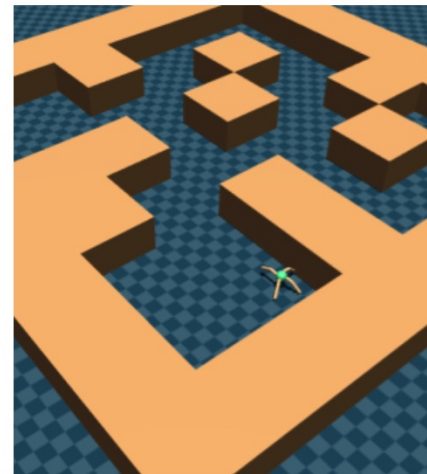
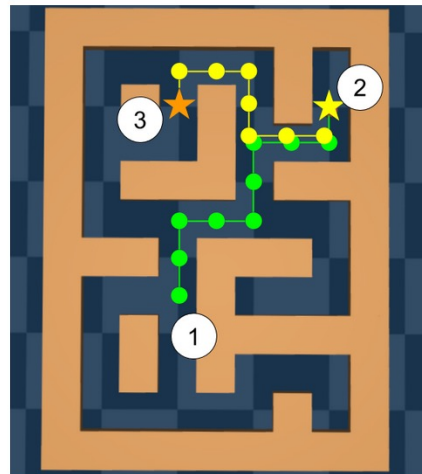


Offline RL?

D4RL

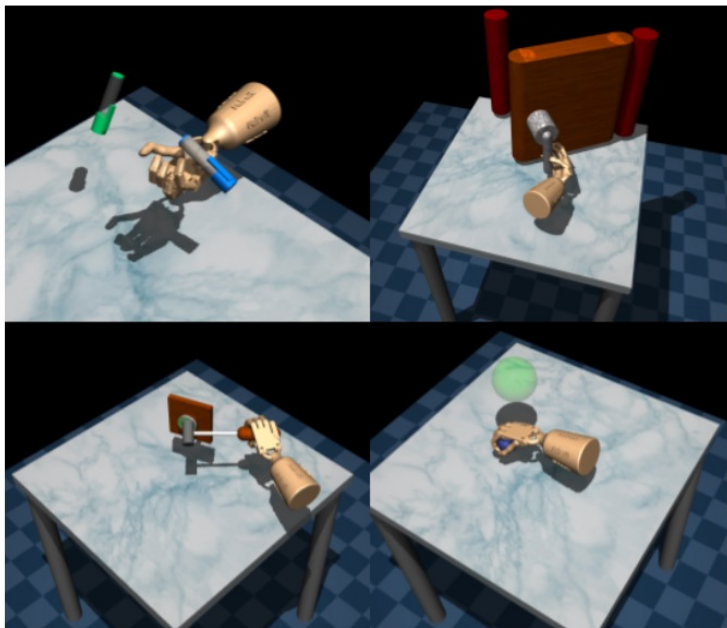
Maze2D and AntMaze

The maze environments are designed to test the ability of agents to recombine existing data in novel ways. For example, if an agent sees trajectories 1-2 and 2-3, it can form a shortest path from 1-3. Two robots are available - a simple ball and the "Ant" robot from the Gym benchmark.



Offline RL?

D4RL

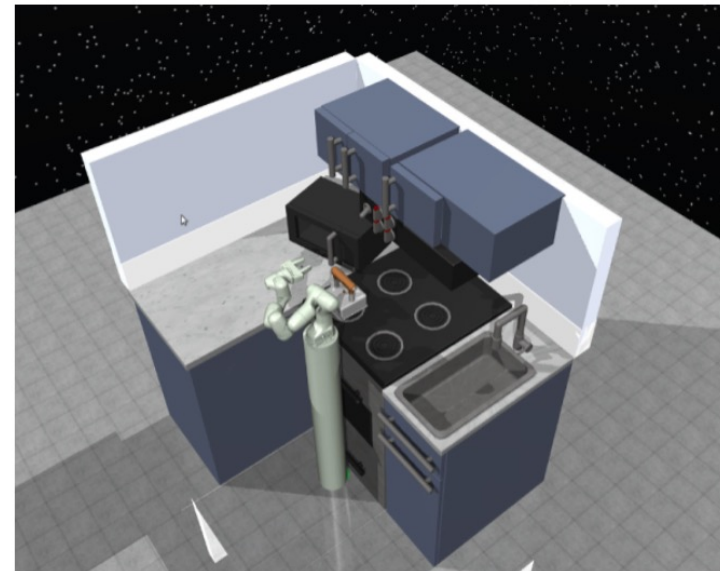


Adroit

The Adroit domain includes motion-captured human data on a realistic, high-DoF robotic hand. A variety of challenging tasks from the original paper are included, including pen twirling, opening a door, using a hammer, and relocating an object.

FrankaKitchen

The FrankaKitchen domain is based on the Adept environment. This domain offers a challenging manipulation problem in an unstructured environment with many possible tasks to perform.



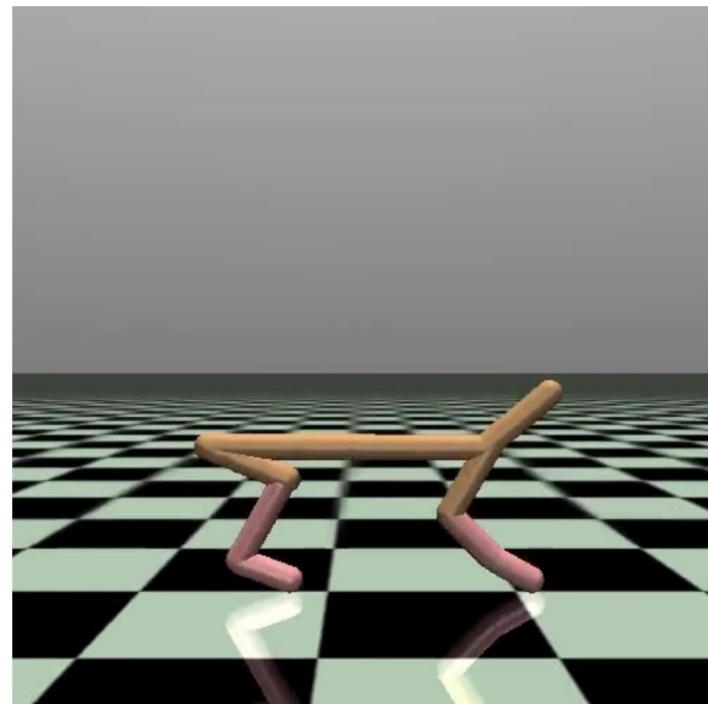


Offline RL?

D4RL

Gym

Several OpenAI Gym benchmark tasks are included with data collected by a variety of pre-trained RL agents. This includes the Hopper, HalfCheetah, and Walker environments.



Why offline RL



Why offline RL



"Place Grapes in Ceramic Bowl"



"Place Bottle In Tray"



"Push Purple Bowl Across The Table"

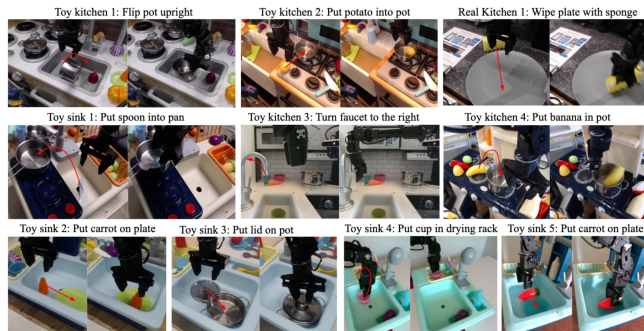


"Wipe Tray With Sponge"

BC-Z dataset

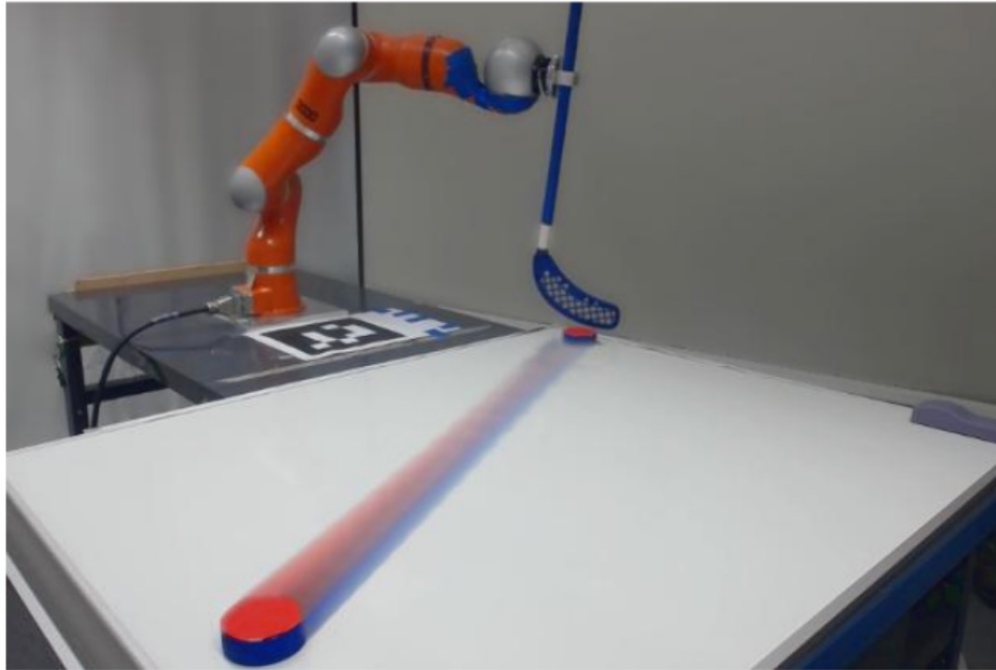


RoboNet



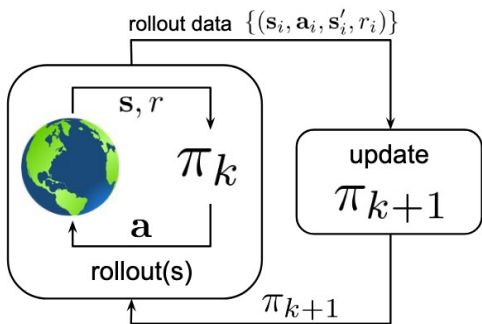
Bridge Dataset

Why offline RL?

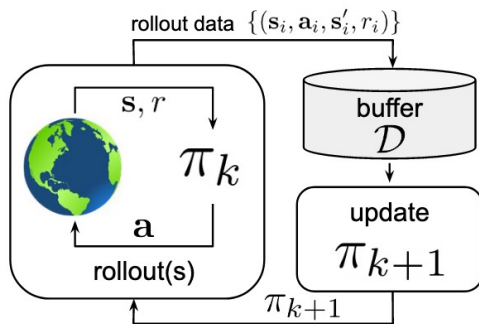


Off-policy Reinforcement Learning

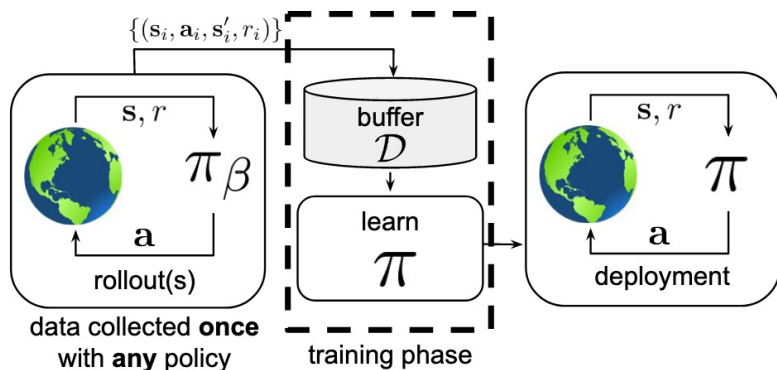
(a) online reinforcement learning



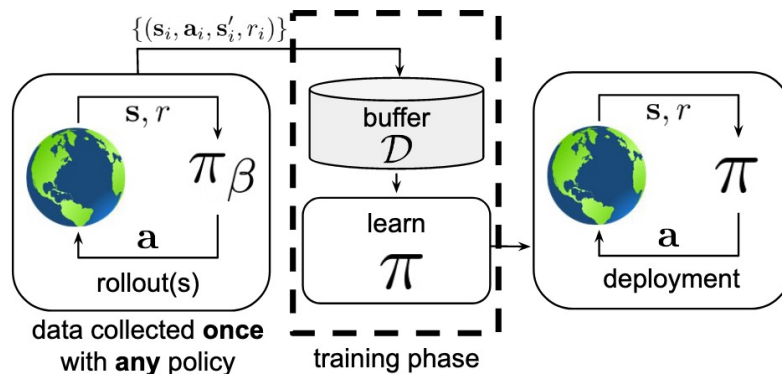
(b) off-policy reinforcement learning



(c) offline reinforcement learning



Actor-Critic RL



$$\mathcal{D} = \{(s, a, s', r)_j\}$$

Critic Update

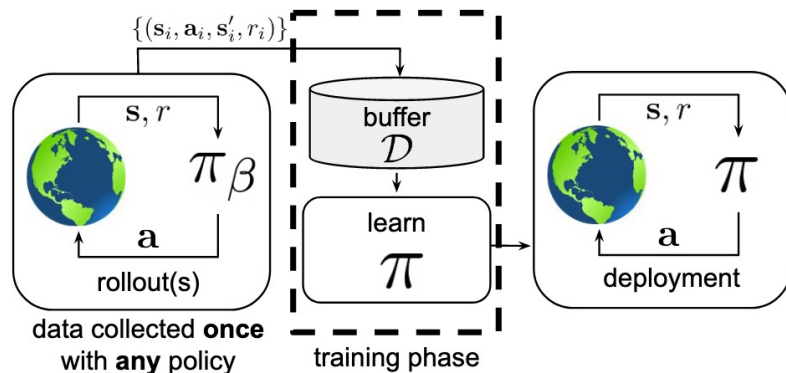
$$\mathcal{T}^\pi Q(s, a) = \mathbb{E}_{s'}[r + \gamma Q(s', \pi(s'))]$$

Critic Update

$$\phi \leftarrow \operatorname{argmax}_\phi \mathbb{E}_{s \in \mathcal{B}}[Q_\theta(s, \pi_\phi(s))]$$

$$Q^\pi(s_t, \mathbf{a}_t) = \mathbb{E}_{\tau \sim p_\pi(\tau | s_t, \mathbf{a}_t)} \left[\sum_{t'=t}^H \gamma^{t'-t} r(s_{t'}, \mathbf{a}_{t'}) \right]$$

Policy Constraints Methods

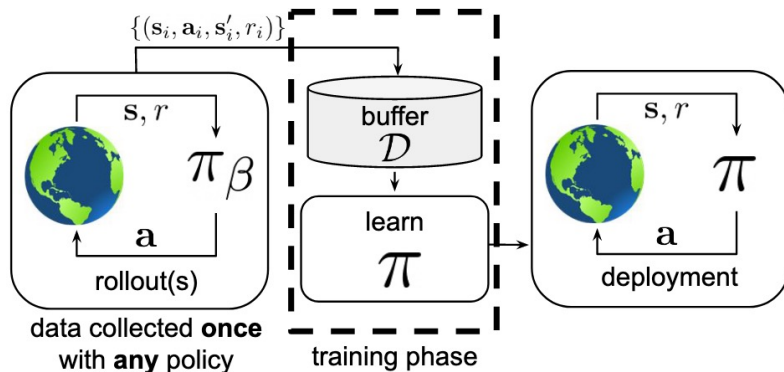


Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems

Sergey Levine^{1,2}, Aviral Kumar¹, George Tucker², Justin Fu¹
¹UC Berkeley, ²Google Research, Brain Team

$$\hat{Q}_{k+1}^\pi \leftarrow \arg \min_Q \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[\left(Q(s, a) - \left(r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_k(a'|s')} [\hat{Q}_k^\pi(s', a')] \right) \right)^2 \right]$$

$$\pi_{k+1} \leftarrow \arg \max_\pi \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{a \sim \pi(a|s)} [\hat{Q}_{k+1}^\pi(s, a)] \right] \text{ s.t. } D(\pi, \pi_\beta) \leq \epsilon.$$



Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems

Sergey Levine^{1,2}, Aviral Kumar¹, George Tucker², Justin Fu¹
¹UC Berkeley, ²Google Research, Brain Team

$$\hat{Q}_{k+1}^\pi \leftarrow \arg \min_Q$$

$$\mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[\left(Q(s, a) - \left(r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_k(a' | s')} [\hat{Q}_k^\pi(s', a')] - \alpha \gamma D(\pi_k(\cdot | s'), \pi_\beta(\cdot | s')) \right) \right)^2 \right]$$

$$\pi_{k+1} \leftarrow \arg \max_{\pi} \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{a \sim \pi(a | s)} [\hat{Q}_{k+1}^\pi(s, a)] - \alpha D(\pi(\cdot | s), \pi_\beta(\cdot | s)) \right].$$



Advantage Weighted Actor Critic

$$\begin{aligned}\pi_{k+1} = \arg \max_{\pi \in \Pi} \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] \\ \text{s.t. } D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_{\beta}(\cdot|\mathbf{s})) \leq \epsilon \\ \int_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) d\mathbf{a} = 1.\end{aligned}$$

$$\begin{aligned}V^{\pi}(\mathbf{s}_t) &= \mathbb{E}_{\tau \sim p_{\pi}(\tau|\mathbf{s}_t)} \left[\sum_{t'=t}^H \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \\ Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) &= \mathbb{E}_{\tau \sim p_{\pi}(\tau|\mathbf{s}_t, \mathbf{a}_t)} \left[\sum_{t'=t}^H \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right]. \quad A(s, a) = Q_{\phi}(s, a) - V(s)\end{aligned}$$



Advantage Weighted Actor Critic

$$\begin{aligned}\mathcal{L}(\pi, \lambda, \alpha) = & \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] \\ & + \lambda(\epsilon - D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_{\beta}(\cdot|\mathbf{s}))) \\ & + \alpha(1 - \int_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) d\mathbf{a}).\end{aligned}$$

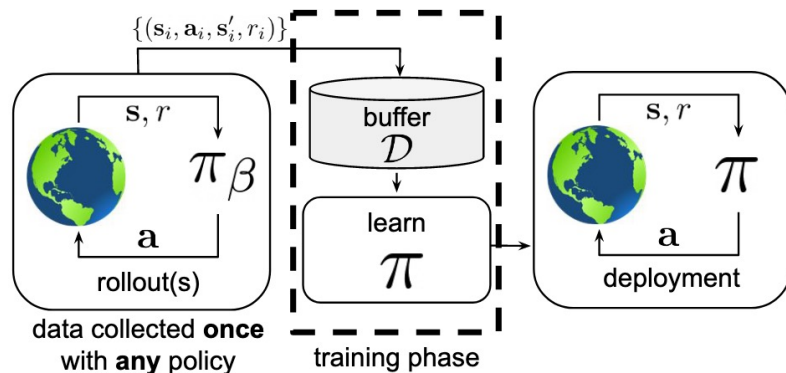
$$\frac{\partial \mathcal{L}}{\partial \pi} = A^{\pi_k}(\mathbf{s}, \mathbf{a}) - \lambda \log \pi_{\beta}(\mathbf{a}|\mathbf{s}) + \lambda \log \pi(\mathbf{a}|\mathbf{s}) + \lambda - \alpha.$$

Advantage Weighted Actor Critic

$$\pi^*(\mathbf{a}|\mathbf{s}) = \frac{1}{Z(\mathbf{s})} \pi_\beta(\mathbf{a}|\mathbf{s}) \exp \left(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a}) \right)$$

$$\begin{aligned} & \arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} [D_{\text{KL}}(\pi^*(\cdot|\mathbf{s}) || \pi_\theta(\cdot|\mathbf{s}))] \\ &= \arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} \left[\mathbb{E}_{\pi^*(\cdot|\mathbf{s})} [-\log \pi_\theta(\cdot|\mathbf{s})] \right] \end{aligned}$$

Policy Constraints Methods



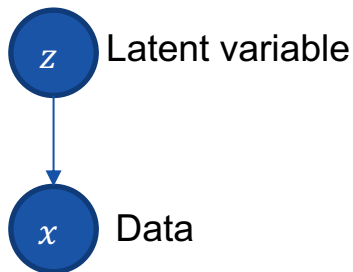
Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems

Sergey Levine^{1,2}, Aviral Kumar¹, George Tucker², Justin Fu¹
¹UC Berkeley, ²Google Research, Brain Team

$$\hat{Q}_{k+1}^\pi \leftarrow \arg \min_Q \mathbb{E}_{(s, a, s') \sim \mathcal{D}} \left[\left(Q(s, a) - \left(r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_k(a'|s')} [\hat{Q}_k^\pi(s', a')] \right) \right)^2 \right]$$

$$\pi_{k+1} \leftarrow \arg \max_\pi \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{a \sim \pi(a|s)} [\hat{Q}_{k+1}^\pi(s, a)] \right] \text{ s.t. } D(\pi, \pi_\beta) \leq \epsilon.$$

The variational lower bound

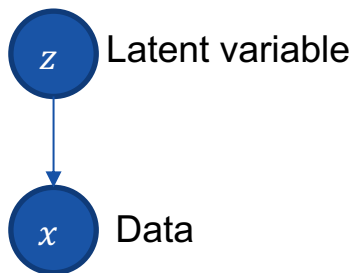


Maximize the loglikelihood of the data

$$\begin{aligned}
 \log p(x) &= \log \int p(x, z) dz \\
 &= \log \int p(x, z) \frac{q(z)}{q(z)} dz \\
 &= \log \int p(x|z)p(z) \frac{q(z|x)}{q(z|x)} dz
 \end{aligned}$$

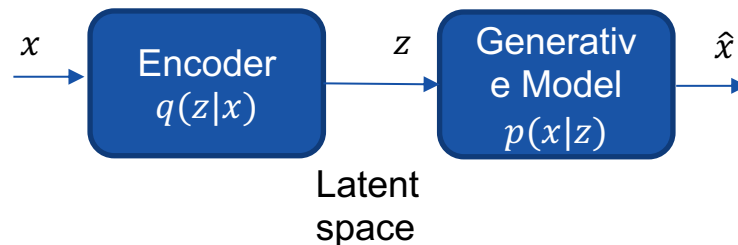
Jensen's inequality

$$\begin{aligned}
 &\geq \mathbb{E}_{q(z|x)}[\log p(x|z)] - \mathbb{E}_{q(z)}\left[\log \frac{q(z)}{p(z)|x}\right] \\
 &= \boxed{\mathbb{E}_{q(z|x)}[\log p(x|z)] - D_{KL}(q(z|x)||p(z))}
 \end{aligned}$$

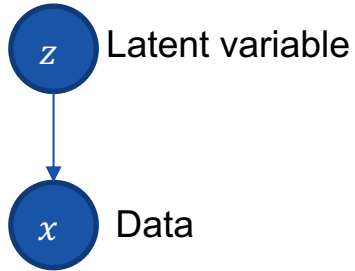


Maximize the variational lower bound

$$= \mathbb{E}_{q(z|x)} [\log p(x|z)] - D_{KL}(q(z|x) || p(z))$$

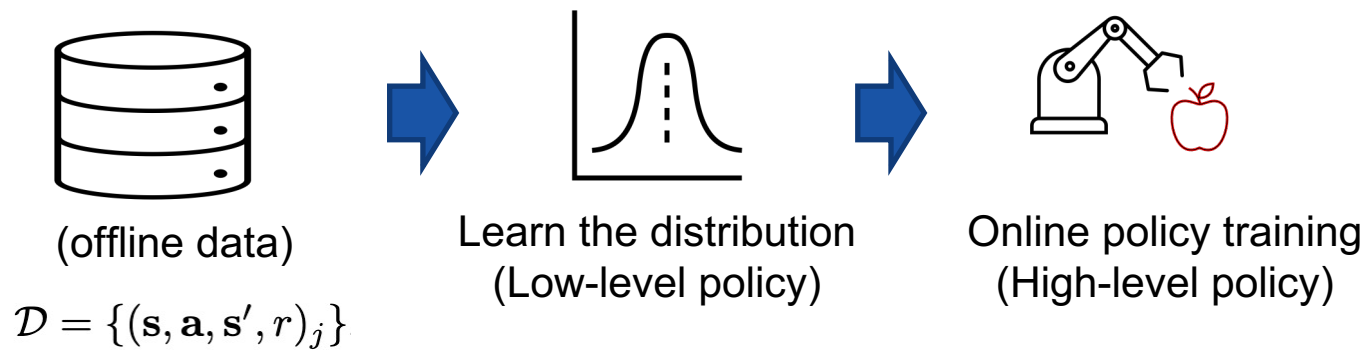


Variational Autoencoder Networks



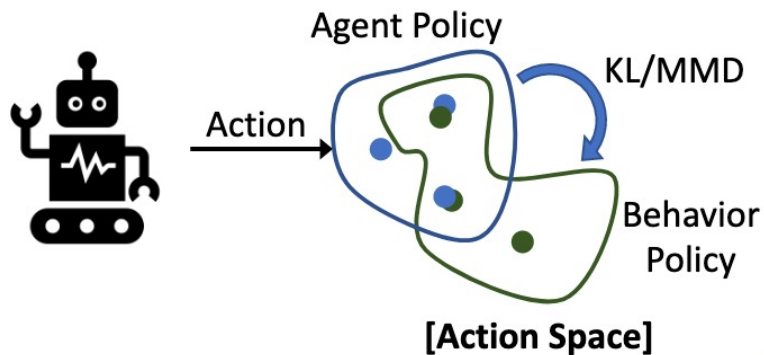


Offline RL with Generative Models

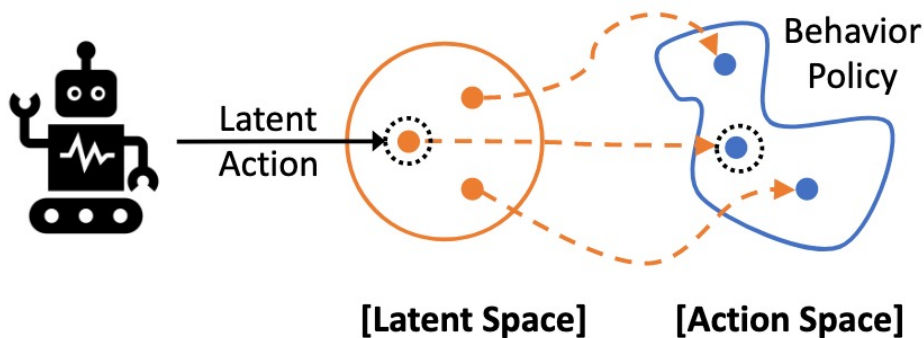


Policy in the Latent Action Space (PLAS)

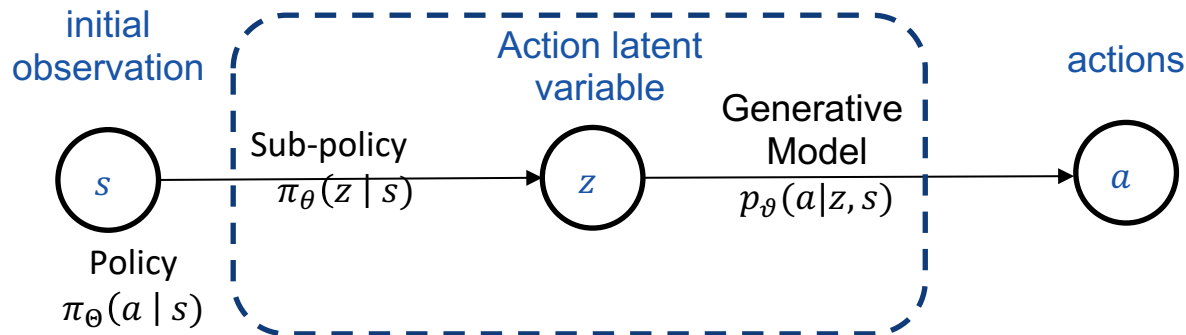
Explicit Policy Constraint



Implicit Policy Constraint using Generative Models



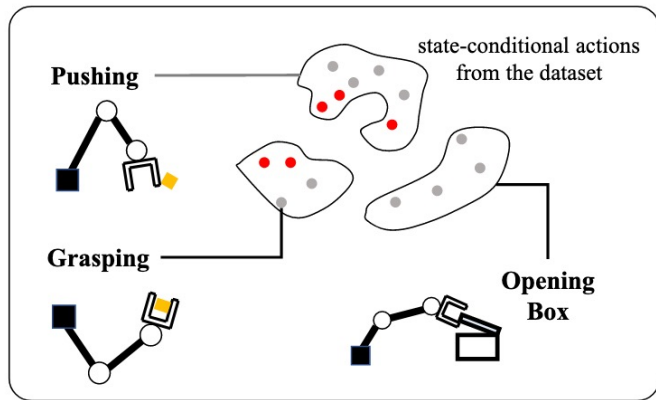
Policy in the Latent Action Space (PLAS)



$$\nabla_{\vartheta} J(\vartheta) = \mathbb{E}_{s \sim \mathcal{D}} [\nabla_a Q(a, s) |_{a=\pi_{\theta}(s, z)} \nabla_z \pi_{\theta}(s, z) |_{z=\pi_{\vartheta}(s)} \nabla_{\vartheta} \pi_{\vartheta}(s)]$$

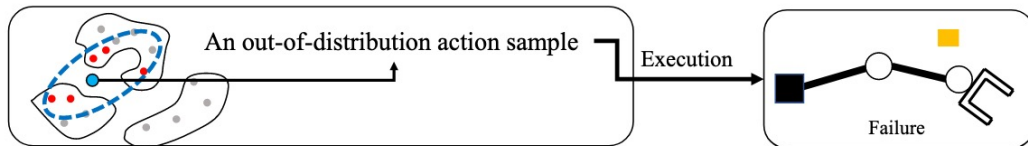
Latent-Variable Advantage-Weighted Policy Optimization

Learning from heterogeneous datasets

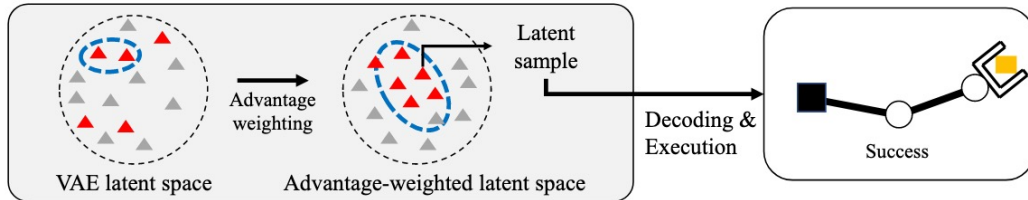


- Task-irrelevant actions
- Task-relevant actions

Action Space Policy Learning



Latent-Variable Advantage-Weighted Policy Optimization (LAPO)



- Prior distribution
- Posterior distribution
- ▲ Task-irrelevant latent actions
- ▲ Task-relevant latent actions

Target Task: Object relocation

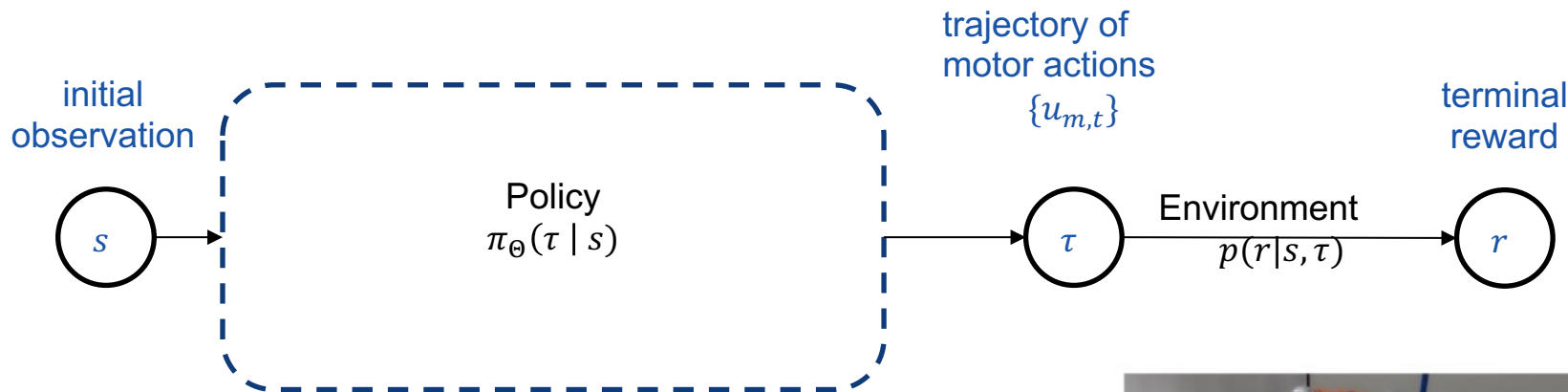
$$\pi^*(a|s) \propto \pi_\beta(a|s) \exp(A(s, a)/\lambda)$$

$$\omega = \exp(A(s, a)/\lambda)$$

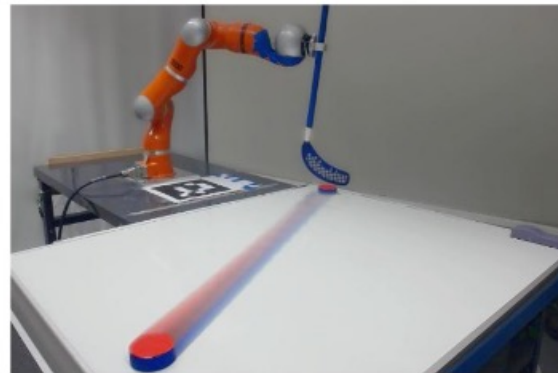
$$\max_{\pi_\theta, q_\psi} \mathbb{E}_{s, a \sim \mathcal{D}} [\omega \mathbb{E}_{q_\psi(z|s, a)} [\log(\pi_\theta(a|s, z)) -$$

$$\beta D_{\text{KL}}(q_\psi(z|s, a) || p(z))]$$

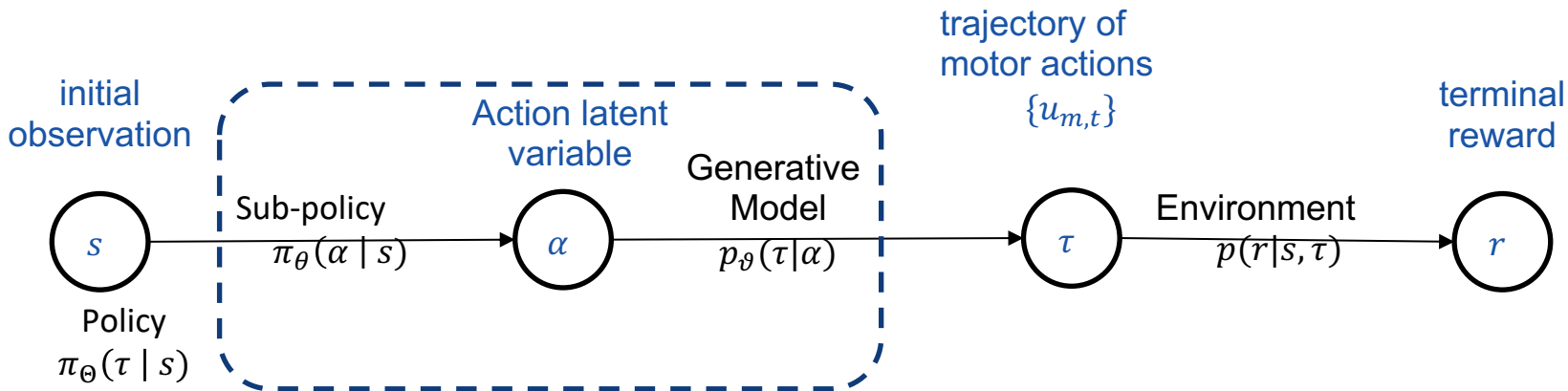
RL with Generative Models



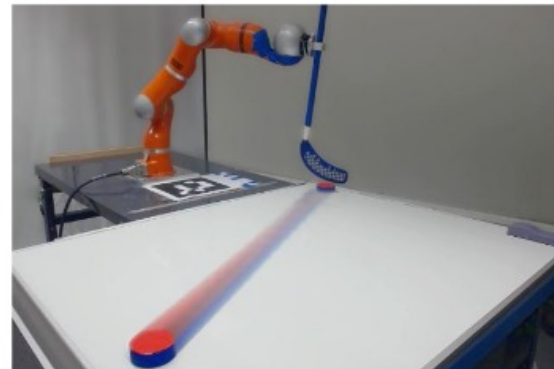
$$\log p(r|s, \Theta) = \log \int p(r|s, \tau) \pi_{\Theta}(\tau|s) d\tau$$



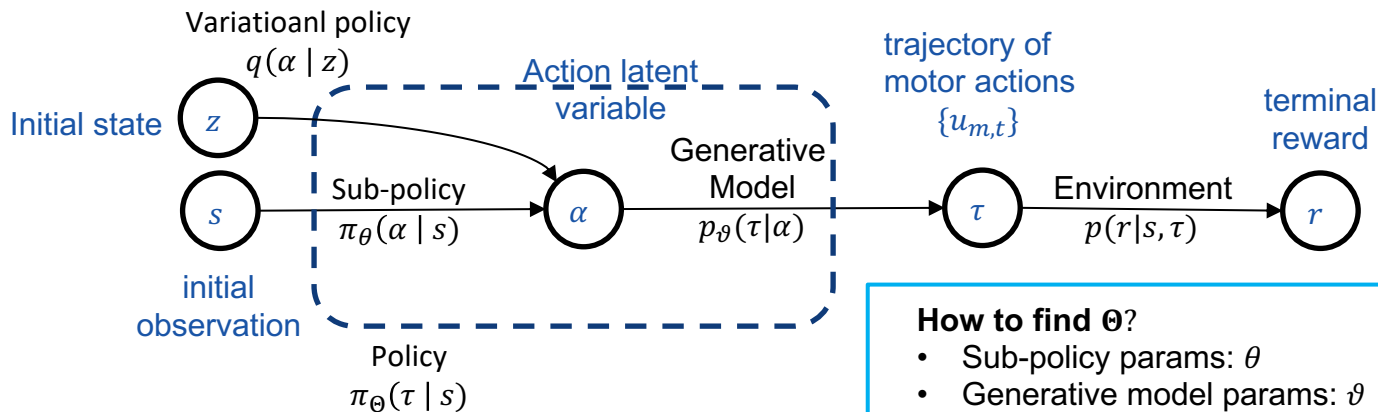
RL with Generative Models



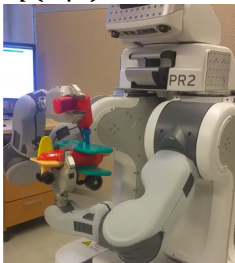
$$\log p(r | s, \Theta) = \log \int p(r | s, p_{\vartheta}(\tau | \alpha)) \pi_{\theta}(\alpha | s) d\alpha$$



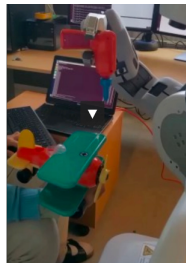
RL with Generative Models



$q(\alpha | z)$



$\pi_{\theta}(\alpha | s)$

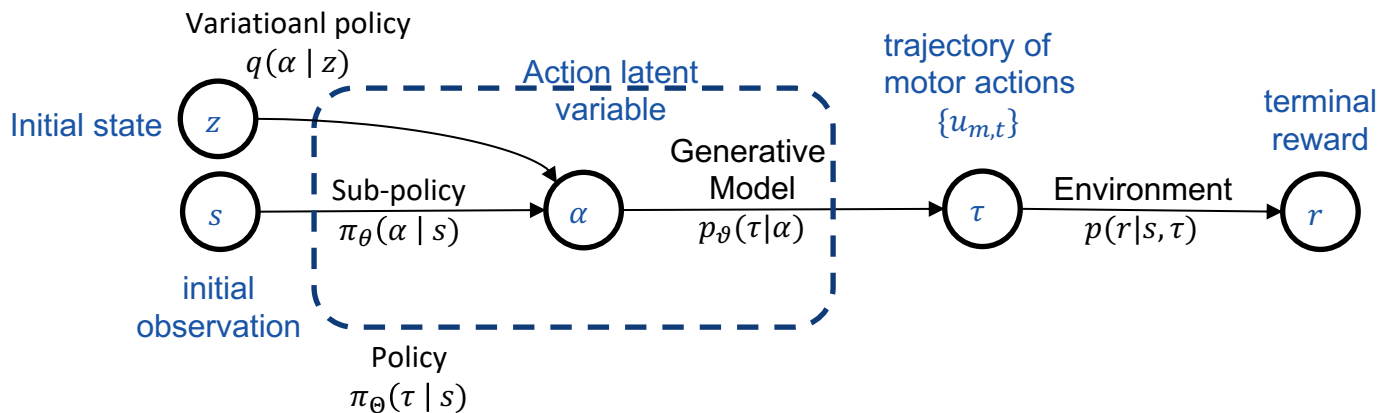


Input
remapping
trick

$$\log p(r | s, \theta) = \log p(r | s, \theta) \int \frac{\log q(\alpha | z)}{\log q(\alpha | z)} q(\alpha | z) d\alpha$$

Ghadirzadeh et al.,
Data-efficient visuomotor policy training
using reinforcement learning and generative models.

RL with Generative Models



$$\log p(r | s, \theta) = \underbrace{\int q(\alpha | z) \log \frac{p(r, \alpha | s, \theta)}{q(\alpha | z)} d\alpha}_{\text{① Lower bound}} + \underbrace{\int q(\alpha | z) \log \frac{q(\alpha | z)}{p(\alpha | r, s, \theta)} d\alpha}_{\text{② KL divergence (non-negative)}}$$

Maximize reward log-likelihood iteratively in two steps:

1. Maximize the lower-bound by minimizing the KL divergence term (update $q(\alpha | z)$)
2. Maximize the lower bound directly (update θ)

RL with Generative Models

- ① Lower bound ② KL divergence (non-negative)

$$\log p(r|s, \theta) = \int q(\alpha|z) \log \frac{p(r, \alpha|s, \theta)}{q(\alpha|z)} d\alpha + \int q(\alpha|z) \log \frac{q(\alpha|z)}{p(\alpha|r, s, \theta)} d\alpha$$

- Expectation step

- Minimize the KL-divergence term by optimizing the variational policy $q(\alpha|z)$

$$q = \operatorname{argmin}_{q'} \int q'(\alpha|z) \log \frac{q'(\alpha|z)}{\pi_{\theta}(\alpha|s)} d\alpha - \int q'(\alpha|z) \log p(r|\alpha, s) d\alpha + \log p(r|s, \theta) \int q'(\alpha|z) d\alpha$$

$$q = \operatorname{argmin}_{q'} D_{KL}(q'(\alpha|z) || \pi_{\theta}(\alpha|s)) - \mathbb{E}_{q'(\alpha|z)}[\log p(r|\alpha, s)]$$

Trust region

Reward seeking

RL with Generative Models

① Lower bound

② KL divergence (non-negative)

$$\log p(r|s, \theta) = \int q(\alpha|z) \log \frac{p(r, \alpha|s, \theta)}{q(\alpha|z)} d\alpha + \int q(\alpha|z) \log \frac{q(\alpha|z)}{p(\alpha|r, s, \theta)} d\alpha$$

- Maximization step

- Maximize the lower bound directly by updating the policy parameters θ

$$\theta = \operatorname{argmax}_{\theta'} \int q(\alpha|z) \log \frac{p(r, \alpha|s, \theta')}{q(\alpha|z)} d\alpha$$

$$= \operatorname{argmax}_{\theta'} \int q(\alpha|z) \log \frac{p(r|\alpha, s) \pi_{\theta'}(\alpha|s)}{q(\alpha|z)} d\alpha$$

$$\theta = \operatorname{argmin}_{\theta'} D_{KL}(q(\alpha|z) || \pi_{\theta'}(\alpha|s))$$

Supervised Learning

$$= \operatorname{argmax}_{\theta'} \int q(\alpha|z) \log \frac{\pi_{\theta'}(\alpha|s)}{q(\alpha|z)} d\alpha + \int q(\alpha|z) \log p(r|\alpha, s) d\alpha$$

RL with Generative Models

① Lower bound

② KL divergence (non-negative)

$$\log p(r|s, \theta) = \int q(\alpha|z) \log \frac{p(r, \alpha|s, \theta)}{q(\alpha|z)} d\alpha + \int q(\alpha|z) \log \frac{q(\alpha|z)}{p(\alpha|r, s, \theta)} d\alpha$$

- Expectation step

$$q = \underset{q'}{\operatorname{argmin}} D_{KL}(q'(\alpha|z) \parallel \pi_{\theta}(\alpha|s)) - \mathbb{E}_{q'(\alpha|z)}[\log p(r|\alpha, s)]$$

Trust region

Reward seeking

- Maximization step

$$\theta = \underset{\theta'}{\operatorname{argmin}} D_{KL}(q(\alpha|z) \parallel \pi_{\theta'}(\alpha|s))$$

Supervised Learning

RL with Generative Models

