

Table of Contents

- 1. Motivation of understanding control as probabilistic inference framework
- 2. Theory of control as probabilistic inference
- 3. Applications, connection to various algorithm
- 4. Future plan for research

Motivation of understanding control as probabilistic inference framework

Maximum Entropy RL

$$J_{\text{MaxEnt}}(\pi; p, r) \triangleq \mathbb{E}_{\mathbf{a_t} \sim \pi(\mathbf{a_t} \mid \mathbf{s_t}), \mathbf{s_{t+1}} \sim p(\mathbf{s_{t+1}} \mid \mathbf{s_t}, \mathbf{a_t})} \left[\sum_{t=1}^{T} r\left(\mathbf{s_t}, \mathbf{a_t}\right) + \alpha \mathcal{H}_{\pi}\left[\mathbf{a_t} \mid \mathbf{s_t}\right] \right],$$

- Representative algorithms: Soft-Q Learning(**SQL**), Soft Actor Critic(**SAC**)
- **KL-Divergence Constraints for Policy Search**

$$\max_{\theta} \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}, a \sim \pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A^{\pi_{\theta_{\text{old}}}}(s, a) \right] \qquad \mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} \left[D_{\text{KL}} \left(\pi_{\theta_{\text{old}}}(\cdot|s) \parallel \pi_{\theta}(\cdot|s) \right) \right] \leq \delta$$

- Representative algorithms: Trust Region Policy Optimization(TRPO), Maximum a posteriori policy optimization (MPO), Advantage weighted regression (AWR),
- **Return conditioned trajectory planning**

$$p\left(\tau \mid \mathcal{O}_{1:T}\right) \propto p(\tau) \exp\left(\sum_{t=1}^{T} r\left(s_{t}, a_{t}\right)\right)$$

Representative algorithms: Planning with diffusion for flexible behavior synthesis (**Diffuser**)

Standard RL Policy Search Problem

Find θ that maximizes cumulative rewards, with parameterization: $p(a_t|s_t,\theta) = \pi_{\theta}(a_t|s_t)$

$$\theta^* = \arg\max_{\theta} \sum_{t=1}^{T} \mathbb{E}_{(s_t, a_t) \sim p(s_t, a_t | \theta)} \left[r(s_t, a_t) \right]$$

In other words,

$$p(\tau|\theta) = p(s_1, a_t, \dots, s_T, a_T \mid \theta) = p(s_1) \prod_{t=1}^T p(a_t \mid s_t, \theta) p(s_{t+1} \mid s_t, a_t).$$
$$\theta^* = \arg\max_{\theta} \mathbb{E}_{\tau \sim p(\tau|\theta)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

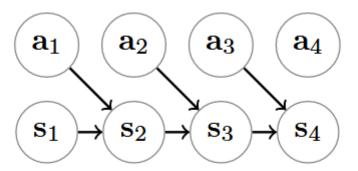
To make RL policy search problem into probabilistic inference framework:

Draw a **Probabilistic Graphical Model (PGM)**, where:

Most probable trajectory: Trajectory from the **optimal policy**.

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the **optimal policy**.

How to draw the PGM?



(a) graphical model with states and actions

- Simplest graphical model to express control problem.
- When we see the graphical model: **Decompose the Joint Distribution!**

$$p(s_1, a_1, s_2, a_2, \dots, a_T, s_{T+1}) = p(s_1) \prod_{t=1}^T p(a_t) \cdot p(s_{t+1} \mid s_t, a_t)$$

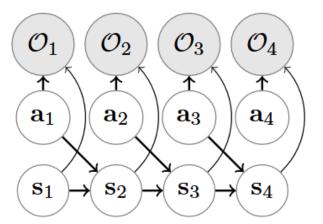
Is it enough to satisfy below questions?

Most probable trajectory: Trajectory from the **optimal policy**.

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the **optimal policy.**

Cannot induce reward/cost information here.

How to draw the PGM?



(b) graphical model with optimality variables

- Introduced additional random variable $\mathcal{O}_t = \begin{cases} 1 & \text{if timestep } t \text{ is optimal} \\ 0 & \text{otherwise} \end{cases}$
- When we see the graphical model: **Decompose the Joint Distribution!**

$$p(\tau, \mathcal{O}_{1:T}) = p(s_1) \prod_{t=1}^{T} p(a_t) p(\mathcal{O}_t = 1 \mid s_t, a_t) p(s_{t+1} \mid s_t, a_t)$$

Decision Choice 1: As we were not interested in action prior (online), we can take it as simplest:

$$p(a_t) = p(a_t|s_t) = \frac{1}{|\mathcal{A}_t|}$$

Decision Choice 2: We define conditional distribution: $p(\mathcal{O}_t = 1 | s_t, a_t) = \exp(r(s_t, a_t))$

Most probable trajectory as trajectory from the optimal policy

From the two Decision Choices, we can factorize posterior:

Decision 1
$$\longrightarrow p(\tau \mid \mathcal{O}_{1:T}) \propto p(\tau, \mathcal{O}_{1:T}) = p(s_1) \prod_{t=1}^{T} p(\mathcal{O}_t = 1 \mid s_t, a_t) p(s_{t+1} \mid s_t, a_t)$$

$$= p(s_1) \prod_{t=1}^{T} \exp(r(s_t, a_t)) p(s_{t+1} \mid s_t, a_t)$$

$$= \left[p(s_1) \prod_{t=1}^{T} p(s_{t+1} \mid s_t, a_t) \right] \exp\left(\sum_{t=1}^{T} r(s_t, a_t) \right).$$

RHS term should be \propto : We can incorporate this defining reward $r(s_t, a_t)$

So we can set the target distribution as:

Product distribution of **Dynamics** and **exponential of Reward Sum (stochastic dynamics)**

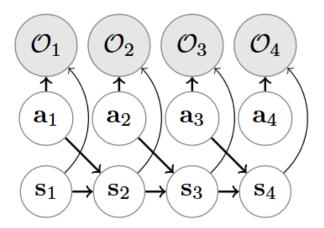
$$p(\tau|o_{1:T}) \propto \left[p(s_1) \prod_{t=1}^T p(s_{t+1} \mid s_t, a_t)\right] \exp\left(\sum_{t=1}^T r(s_t, a_t)\right).$$

Product distribution of Feasible dynamics and exponential of Reward Sum (deterministic dynamics)

$$p(\tau|o_{1:T}) \propto \mathbb{1}[p(\tau) \neq 0] \exp\left(\sum_{t=1}^{T} r(s_t, a_t)\right)$$

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the optimal policy.

Now lets find out how to recover **optimal policy** $\mathbf{p}(\mathbf{a_t}|\mathbf{s_t}, \mathcal{O}_{\mathbf{t}:\mathbf{T}} = \mathbf{1})$ from designed PGM.



(b) graphical model with optimality variables

- Definition 1. $\beta_t(s_t, a_t) = p(\mathcal{O}_{t:T}|s_t, a_t)$
- Definition 2. $\beta_t(s_t) = p(\mathcal{O}_{t:T}|s_t)$

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the optimal policy.

Definition 1. $\beta_t(s_t, a_t) = p(\mathcal{O}_{t:T}|s_t, a_t)$ Definition 2. $\beta_t(s_t) = p(\mathcal{O}_{t:T}|s_t)$

Formulation 1; Marginalizing $\beta_t(s_t)$ with a_t

$$\beta_t(s_t) = p(\mathcal{O}_{t:T} \mid s_t) = \int_{\mathcal{A}} p(\mathcal{O}_{t:T}, a_t | s_t) da_t \tag{1}$$

$$= \int_{\mathcal{A}} \frac{p(\mathcal{O}_{t:T}, a_t|s_t)}{p(a_t|s_t)} p(a_t|s_t) da_t \tag{2}$$

$$= \int_{\Lambda} p(\mathcal{O}_{t:T} \mid s_t, a_t) p(a_t \mid s_t) \, da_t \tag{3}$$

$$= \int_{\mathcal{A}} \beta_t(s_t, a_t) p(a_t \mid s_t) \, da_t. \tag{4}$$

Decision 1
$$\Longrightarrow$$
 $\beta_t(s_t) = \int_{\mathcal{A}} \beta_t(s_t, a_t) da_t$

$$p(a_t) = p(a_t|s_t) = \frac{1}{|\mathcal{A}_t|}$$

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the optimal policy.

Definition 1.
$$\beta_t(s_t, a_t) = p(\mathcal{O}_{t:T}|s_t, a_t)$$
 Definition 2. $\beta_t(s_t) = p(\mathcal{O}_{t:T}|s_t)$

• **Formulation 2**; Marginalizing $\beta_t(s_t, a_t)$ with s_{t+1}

$$\beta_t(s_t, a_t) = p(\mathcal{O}_{t:T} \mid s_t, a_t) \tag{1}$$

$$= \int_{\mathcal{S}} \beta_{t+1}(s_{t+1}) \, p(s_{t+1} \mid s_t, a_t) \, p(\mathcal{O}_t \mid s_t, a_t) \, ds_{t+1}. \tag{2}$$

• Formulation 3; Optimal target policy from $\beta_t(s_t, a_t)$, $\beta_t(s_t)$

$$p(a_t \mid s_t, \mathcal{O}_{t:T}) = \frac{p(s_t, a_t \mid \mathcal{O}_{t:T})}{p(s_t \mid \mathcal{O}_{t:T})} = \frac{p(\mathcal{O}_{t:T} \mid s_t, a_t) p(a_t \mid s_t) p(s_t)}{p(\mathcal{O}_{t:T} \mid s_t, a_t)}$$

$$\propto \frac{p(\mathcal{O}_{t:T} \mid s_t, a_t)}{p(\mathcal{O}_{t:T} \mid s_t)} = \frac{\beta_t(s_t, a_t)}{\beta_t(s_t)},$$

Decision 1

Inference of posterior action conditional $p(a_t|s_t,\theta)$: Gives the optimal policy.

- Definition 3 (Soft Q): $Q(\mathbf{s}_t, \mathbf{a}_t) = \log \beta_t(\mathbf{s}_t, \mathbf{a}_t)$
- Definition 4 (Soft V): $V(\mathbf{s}_t) = \log \beta_t(\mathbf{s}_t)$
- From formulation 1:

$$\beta_t(s_t) = \int_{\mathcal{A}} \beta_t(s_t, a_t) da_t$$

If the scale of Q is large, High Q-value dominates the other actions.



Soft Maximization

$$V(\mathbf{s}_t) = \log \int_{\mathcal{A}} \exp(Q(\mathbf{s}_t, \mathbf{a}_t)) d\mathbf{a}_t \approx \max_{\mathbf{a}_t} Q(\mathbf{s}_t, \mathbf{a}_t).$$

From formulation 2:

$$\beta_t(s_t, a_t) = \int_{\mathcal{S}} \beta_{t+1}(s_{t+1}) p(s_{t+1} \mid s_t, a_t) p(\mathcal{O}_t \mid s_t, a_t) ds_{t+1}.$$

$$Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \log E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)} \left[\exp \left(V(\mathbf{s}_{t+1}) \right) \right]$$

From formulation 3:

$$p(a_t \mid s_t, \mathcal{O}_{t:T}) = \frac{\beta_t(s_t, a_t)}{\beta_t(s_t)} = \exp(Q(s_t, a_t) - V(s_t))$$

Problem of the drawn PGM; Overly optimistic Q function

Stochastic dynamics Leads to overly optimistic Q function, risk seeking property.

$$V\left(\mathbf{s}_{t}\right) = \log \int_{\mathcal{A}} \exp\left(Q\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right)\right) d\mathbf{a}_{t} \approx \max_{\mathbf{a}_{t}} Q\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right).$$
 Another soft maximization term
$$Q\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) = r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \log E_{\mathbf{s}_{t+1} \sim p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}\right)} \left[\exp\left(V\left(\mathbf{s}_{t+1}\right)\right)\right]$$

Where original bellman expectation equation:

$$V(s_t) = \max_{a_t} Q(s_t, a_t)$$

$$Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + E_{\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)} [V(\mathbf{s}_{t+1})]$$

- EX): Two states for s_{t+1} ,
 - Win the lottery with $p(s_{t+1} = win | s_t, a_t = buy) = 0.000001$
 - Lose the lottery with $p(s_{t+1} = lose \mid s_t, a_t = buy) = 0.9999999$

As $V(s_{t+1} = win)$ dominates, always choose action to buy lottery.

Problem of the drawn PGM; Overly optimistic Q function

Deterministic dynamics much more desirable.

$$V(\mathbf{s}_{t}) = \log \int_{\mathcal{A}} \exp \left(Q(\mathbf{s}_{t}, \mathbf{a}_{t})\right) d\mathbf{a}_{t} \approx \max_{\mathbf{a}_{t}} Q(\mathbf{s}_{t}, \mathbf{a}_{t}).$$
$$Q(\mathbf{s}_{t}, \mathbf{a}_{t}) = r(\mathbf{s}_{t}, \mathbf{a}_{t}) + V(\mathbf{s}_{t+1})$$

Where original bellman expectation equation:

$$V(s_t) = \max_{a_t} Q(s_t, a_t)$$
$$Q(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + V(\mathbf{s}_{t+1})$$

We can **safely** use deterministic update of bellman equation, extract the policy.

$$p(a_t \mid s_t, \mathcal{O}_{t:T}) = \frac{\beta_t(s_t, a_t)}{\beta_t(s_t)} = \exp(Q(s_t, a_t) - V(s_t))$$

Then, is this framework only available in deterministic settings?

So what this inference of optimal conditioned policy exactly optimizes?

Lets get back to the optimality conditioned trajectory distribution (**Target**).

$$p(\tau|o_{1:T}) = \left[p(\mathbf{s}_1) \prod_{t=1}^{T} p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)\right] \exp\left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)\right)$$
(1)

What we want **to model is policy**; So we formulate differently with above equation.

$$p(\tau \mid \mathcal{O}_{1:T}) = p(s_1 \mid \mathcal{O}_{1:T}) \cdot \prod_{t=1}^{T} p(a_t \mid s_t, \mathcal{O}_{1:T}) \cdot p(s_{t+1} \mid s_t, a_t, \mathcal{O}_{1:T})$$

We can define approximate trajectory distribution closest to the target distribution.

$$\hat{p}(\tau) = p\left(\mathbf{s}_1 \mid \mathcal{O}_{1:T}\right) \prod_{t=1}^{T} p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t, \mathcal{O}_{1:T}\right) \pi\left(\mathbf{a}_t \mid \mathbf{s}_t\right)$$
(2)

Then make this optimization problem: (1), (2)

$$D_{\mathrm{KL}}(\hat{p}(\tau)||p(\tau)) = -E_{\tau \sim \hat{p}(\tau)}[\log p(\tau) - \log \hat{p}(\tau)].$$

So what this inference of optimal conditioned policy exactly optimizes?

In deterministic case,

$$p(\tau|o_{1:T}) = \left[p(\mathbf{s}_1) \prod_{t=1}^{T} p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)\right] \exp\left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)\right)$$
(1)

Also most probable approximate distribution will be:

$$\hat{p}(\tau) = p\left(\mathbf{s}_{1} \mid \mathcal{O}_{1:T}\right) \prod_{t=1}^{T} p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}, \mathcal{O}_{1:T}\right) \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)$$

$$= \hat{p}(\tau) = p\left(\mathbf{s}_{1}\right) \prod_{t=1}^{T} p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}\right) \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)$$
(3)

Then optimization problem solves: (1), (3)

$$D_{\mathrm{KL}}(\hat{p}(\tau)||p(\tau)) = -E_{\tau \sim \hat{p}(\tau)}[\log p(\tau) - \log \hat{p}(\tau)].$$

So what this inference of optimal conditioned policy exactly optimizes?

$$-D_{\mathrm{KL}}(\hat{p}(\tau)||p(\tau)) = E_{\tau \sim \hat{p}(\tau)} \left[\log p\left(\mathbf{s}_{1}\right) + \sum_{t=1}^{T} \left(\log p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}\right) + r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right)\right) - \log p\left(\mathbf{s}_{1}\right) - \sum_{t=1}^{T} \left(\log p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}\right) + \log \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right) \right]$$

$$= E_{\tau \sim \hat{p}(\tau)} \left[\sum_{t=1}^{T} r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) - \log \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right) \right]$$

$$= \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim \hat{p}(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) - \log \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right) \right]$$

$$= \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim \hat{p}(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + E_{\mathbf{s}_{t} \sim \hat{p}(\mathbf{s}_{t})} \left[\mathcal{H}\left(\pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right) \right]$$

In deterministic setting; Approximation of inference, equivalent to max entropy RL

So what this inference of optimal conditioned policy exactly optimizes?

In stochastic case we cannot derive the same as deterministic.

$$p(\tau|o_{1:T}) = \left[p(\mathbf{s}_1) \prod_{t=1}^{T} p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)\right] \exp\left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)\right)$$
(1)

$$\hat{p}(\tau) = p\left(\mathbf{s}_1 \mid \mathcal{O}_{1:T}\right) \prod_{t=1}^{T} p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t, \mathcal{O}_{1:T}\right) \pi\left(\mathbf{a}_t \mid \mathbf{s}_t\right)$$
(2)

- We cannot erase optimality conditioned dynamics terms.
- Directly optimize: $D_{\mathrm{KL}}(\hat{p}(\tau)||p(\tau)) = -E_{\tau \sim \hat{p}(\tau)}[\log p(\tau) \log \hat{p}(\tau)].$
- Leads: $E_{\tau \sim \hat{p}(\tau)} \left[\log p\left(\mathbf{s}_{1}\right) + \sum_{t=1}^{T} r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \log p\left(\mathbf{s}_{t+1} \mid \mathbf{s}_{t}, \mathbf{a}_{t}\right) \right] + \mathcal{H}(\hat{p}(\tau))$

So, to modeling the approximation as (2), we cannot use model-free algorithms.

Direct Solution! Abuse PGM.

Why not we just model our approxmate distribution with (3) rather than (2)?

$$p(\tau|o_{1:T}) = \left[p\left(\mathbf{s}_{1}\right)\prod_{t=1}^{T}p\left(\mathbf{s}_{t+1}\mid\mathbf{s}_{t},\mathbf{a}_{t}\right)\right] \exp\left(\sum_{t=1}^{T}r\left(\mathbf{s}_{t},\mathbf{a}_{t}\right)\right) (1)$$

$$\hat{p}(\tau) = p\left(\mathbf{s}_{1}\mid\mathcal{O}_{1:T}\right)\prod_{t=1}^{T}p\left(\mathbf{s}_{t+1}\mid\mathbf{s}_{t},\mathbf{a}_{t},\mathcal{O}_{1:T}\right)\pi\left(\mathbf{a}_{t}\mid\mathbf{s}_{t}\right) (2)$$

$$= \hat{p}(\tau) = p\left(\mathbf{s}_{1}\right)\prod_{t=1}^{T}p\left(\mathbf{s}_{t+1}\mid\mathbf{s}_{t},\mathbf{a}_{t}\right)\pi\left(\mathbf{a}_{t}\mid\mathbf{s}_{t}\right) (3)$$

- Pros: We can recover the maximum entropy RL objective, like deterministic case. Agents cannot control dynamical systems like (2), not intuitive. (Is there anyone who can make environment to always win the lottery?)
- Cons: We abuse the PGM w.r.t. optimality -> No theoretical grounds (yet)

Connection to Structural Variational Inference

- The theoretical grounds of abused solution!
- Target distribution: $p(\tau|o_{1:T}) = \left[p(\mathbf{s}_1) \prod_{t=1}^{T} p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)\right] \exp\left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)\right)$ (1)
- Variational approximate distribution: $\hat{q}(\tau) = q(\mathbf{s}_1) \prod_{t=1}^{T} q(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t \mid \mathbf{s}_t)$. (2) q, not p
- With given observation, **ELBO with variational trajectory distribution:**

$$\log p\left(\mathcal{O}_{1:T}\right) = \log \iint p\left(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right) d\mathbf{s}_{1:T} d\mathbf{a}_{1:T}$$

$$= \log \iint p\left(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right) \frac{q\left(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right)}{q\left(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right)} d\mathbf{s}_{1:T} d\mathbf{a}_{1:T}$$

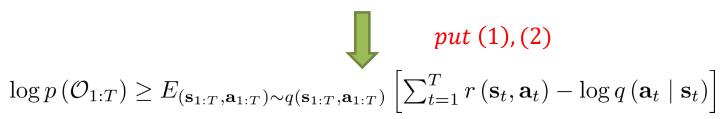
$$= \log E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[\frac{p\left(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right)}{q\left(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right)} \right]$$

$$\geq E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[\log p\left(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right) - \log q\left(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right) \right]$$

Connection to Structural Variational Inference

- The theoretical grounds of abused solution!
- Target distribution: $p(\tau|o_{1:T}) = \left[p(\mathbf{s}_1) \prod_{t=1}^{T} p(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t)\right] \exp\left(\sum_{t=1}^{T} r(\mathbf{s}_t, \mathbf{a}_t)\right)$ (1)
- Variational approximate distribution: $\hat{q}(\tau) = q(\mathbf{s}_1) \prod_{t=1}^{T} q(\mathbf{s}_{t+1} \mid \mathbf{s}_t, \mathbf{a}_t) q(\mathbf{a}_t \mid \mathbf{s}_t)$. (2) q, not p

$$\log p(\mathcal{O}_{1:T}) \ge E_{(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}) \sim q(\mathbf{s}_{1:T}, \mathbf{a}_{1:T})} \left[\log p\left(\mathcal{O}_{1:T}, \mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right) - \log q\left(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}\right)\right]$$



• Then we set dynamics distributions q with $p(s_1), p(s_{t+1}|s_t, a_t) \to \text{Recovers Max-Entropy RL}$

$$= \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim \hat{p}(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \mathcal{H}\left(\pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right) \right]$$

Abused distribution is one of the ELBO, derivation of our designed PGM (Got Theoretical Grounds.)

How can we extract policy from the abused settings (ELBO)

- We already discussed before **soft-Q**, V and optimal policy for **stochastic/deterministic** settings.
- What about **abused settings**, any difference with $\exp(Q(s_t, a_t) V(s_t))$?

From the max-entropy objective function:

$$\sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim \hat{p}(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \mathcal{H}\left(\pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right) \right]$$

If t = T:

$$\mathbb{E}_{(s_T, a_T) \sim \hat{p}(s_T, a_T)} \left[r(s_T, a_T) - \log \pi(a_T \mid s_T) \right]$$

$$= \mathbb{E}_{s_T \sim \hat{p}(s_T)} \left[\mathbb{E}_{a_T \sim \pi(\cdot \mid s_T)} \left[r(s_T, a_T) - \log \pi(a_T \mid s_T) \right] \right]$$

$$= \mathbb{E}_{s_T \sim \hat{p}(s_T)} \left[\mathbb{E}_{a_T \sim \pi(\cdot \mid s_T)} \left[-\log \pi(a_T \mid s_T) + r(s_T, a_T) \right] \right]$$

$$= \mathbb{E}_{s_T \sim \hat{p}(s_T)} \left[-D_{\text{KL}} \left(\pi(a_T \mid s_T) \parallel \frac{1}{Z(s_T)} \exp(r(s_T, a_T)) \right) + \log Z(s_T) \right]$$

$$= \mathbb{E}_{s_T \sim \hat{p}(s_T)} \left[-D_{\text{KL}} \left(\pi(a_T \mid s_T) \parallel \frac{1}{\exp(V(s_T))} \exp(r(s_T, a_T)) \right) + V(s_T) \right]$$

$$\pi^*(a_T \mid s_t) = \exp(r(s_T, a_T))$$

Define:

$$V(s_T) = \log \int_{\mathcal{A}} \exp(r(s_T, a_t)) da_T$$

$$\pi^*(a_T|s_t) = \exp(r(s_T, a_T))$$

How can we extract policy from the abused settings (ELBO)

If t = t (intermediate steps):

• Policy should also consider the **future expected value**. (Summation of $t = 1 \sim T$ recovers original objective.)

$$\mathbb{E}_{(s_t,a_t)\sim \hat{p}(s_t,a_t)}\left[r(s_t,a_t) - \log\pi(a_t\mid s_t) + \mathbb{E}_{s_{t+1}\sim p(\cdot\mid s_t,a_t)}\left[V(s_{t+1})\right]\right]$$

$$\stackrel{(1)}{=} \mathbb{E}_{(s_t,a_t)\sim \hat{p}}\left[\underline{r(s_t,a_t) + \mathbb{E}_{s_{t+1}}[V(s_{t+1})]} - \log\pi(a_t\mid s_t)\right]$$

$$:= Q(s_t,a_t) \qquad \qquad \text{Define, same as Bellman Expectation Eq.}$$

$$\text{We can mitigate risk-seeking Q Problem!}$$

$$\stackrel{(2)}{=} \mathbb{E}_{s_t\sim \hat{p}(s_t)}\left[-D_{\text{KL}}\left(\pi(\cdot\mid s_t) \left\| \frac{1}{Z(s_t)} \exp(Q(s_t,\cdot))\right) + \log Z(s_t)\right]\right] \qquad \text{Now defined on Q}$$

$$= \mathbb{E}_{s_t\sim \hat{p}(s_t)}\left[-D_{\text{KL}}\left(\pi(a_t\mid s_t) \left\| \frac{1}{\exp(V(s_t))} \exp(Q(s_t,a_t))\right) + V(s_t)\right]\right] \qquad \qquad \pi^*\left(\mathbf{a}_t\mid \mathbf{s}_t\right) = \exp\left(Q\left(\mathbf{s}_t,\mathbf{a}_t\right) - V\left(\mathbf{s}_t\right)\right),$$

In Summary:

- 1. We want to transform the control problem -> Probabilistic Inference.
- 2. Draw PGM which can infer $p(a_t|s_t, O_{1:T})$, $p(\tau|O_{1:T})$ posteriors.
- 3. Decompose and **got optimal target distribution**: $p(\tau|O_{1:T})$
- 4. To find out optimal policy, defined **soft** Q, V functions and found **analytical form of** $\pi^*(a_t|s_t)$ Stochastic dynamics has risk-seeking problem with Q function.
- 5. With defining optimal approximate distribution (closest distribution), Transformed inference problem to optimization problem. Again, problem with stochastic dynamics setting has problem. (Only model-based works)
- 6. We abused the PGM, and modeled approximate distribution differently. This reveals that one of the **ELBO** of PGM := Abused distribution.
- 7. Abused distribution recovers Max Entropy RL, no problem of risk-seeking Q function. Most of the Max-Entropy RL algorithms built on this setting.

Maximum Entropy Policy Gradients

Policy Gradients with the Max-Entropy Objective

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^{T} \nabla_{\theta} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \mathcal{H}\left(q_{\theta}\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right) \right]$$

$$= \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[\nabla_{\theta} \log q_{\theta}\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right) \left(\sum_{t'=t}^{T} r\left(\mathbf{s}_{t'}, \mathbf{a}_{t'}\right) - \log q_{\theta}\left(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}\right) - 1 \right) \right]$$

$$= \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[\nabla_{\theta} \log q_{\theta}\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right) \left(\sum_{t'=t}^{T} r\left(\mathbf{s}_{t'}, \mathbf{a}_{t'}\right) - \log q_{\theta}\left(\mathbf{a}_{t'} \mid \mathbf{s}_{t'}\right) - b\left(\mathbf{s}_{t'}\right) \right) \right],$$

$$\nabla_{\theta} J(\theta) = \sum_{t=1}^{T} E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[\nabla_{\theta} \log q_{\theta} \left(\mathbf{a}_{t} \mid \mathbf{s}_{t} \right) \hat{A} \left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) \right]$$

• We can just put $-\log q_{\theta}(a_t|s_t)$ penalty to the original reward, and use the policy gradient algorithms.

Soft Q Learning

• With below two equations, definition of soft V, and its optimal policy $q(a_t|s_t)$.

$$V(\mathbf{s}_t) = \log \int_{\mathcal{A}} \exp(Q(\mathbf{s}_t, \mathbf{a}_t)) d\mathbf{a}_t$$
$$q(\mathbf{a}_t \mid \mathbf{s}_t) = \exp(Q(\mathbf{s}_t, \mathbf{a}_t) - V(\mathbf{s}_t))$$

- We can parameterize Q_{ϕ} only to express both equations.
- Training Q_{ϕ} : $\mathcal{E}(\phi) = E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[\left(r\left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) + E_{q(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})} \left[V_{\psi}\left(\mathbf{s}_{t+1} \right) \right] Q_{\phi}\left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) \right]$ $\phi \leftarrow \phi \alpha E \left[\frac{dQ_{\phi}}{d\phi} \left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) \left(Q_{\phi}\left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) \left(r\left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) + \underline{\log \int_{\mathcal{A}} \exp\left(Q\left(\mathbf{s}_{t+1}, \mathbf{a}_{t+1} \right) \right) d\mathbf{a}_{t+1}} \right) \right) \right].$ $\log \mathbb{E}_{q(a_{t})} \frac{\exp\left(Q\left(\mathbf{s}_{t+1}, \mathbf{a}_{t+1} \right) \right)}{q(a_{t})}$
- Training π_{θ} :

Stein Variational Gradient Descent (**SVGD**) to match: $\pi_{\theta}(a_t|s_t) \propto \exp(Q_{soft}^{\phi}(s_t,\cdot))$

Importance Sampling with current policy

Soft Actor Critic

• Explicit modeling of the policy network $\pi(a_t|s_t;\theta)$

$$\mathcal{T}^{\pi}Q\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) \triangleq r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p}\left[V\left(\mathbf{s}_{t+1}\right)\right]$$

$$V\left(\mathbf{s}_{t}\right) = \mathbb{E}_{\mathbf{a}_{t} \sim \pi}\left[Q\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) - \log \pi\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right)\right] \qquad \text{Check lemma 1 from Haarnoja (2018)}$$

Train both soft V, Q with MSE loss:

$$\mathcal{E}(\phi) = E_{(\mathbf{s}_{t}, \mathbf{a}_{t}) \sim q(\mathbf{s}_{t}, \mathbf{a}_{t})} \left[\left(r\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) + E_{q(\mathbf{s}_{t+1}|\mathbf{s}_{t}, \mathbf{a}_{t})} \left[V_{\psi}\left(\mathbf{s}_{t+1}\right) \right] - Q_{\phi}\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right)^{2} \right]$$

$$\mathcal{E}(\psi) = E_{\mathbf{s}_{t} \sim q(\mathbf{s}_{t})} \left[\left(E_{\mathbf{a}_{t} \sim q(\mathbf{a}_{t}|\mathbf{s}_{t})} \left[Q_{\phi}\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right) - \log q\left(\mathbf{a}_{t} \mid \mathbf{s}_{t}\right) \right] - V_{\psi}\left(\mathbf{s}_{t}, \mathbf{a}_{t}\right)^{2} \right].$$

• Training π_{θ} :

$$J_{\pi}(\theta) = \mathbb{E}_{\mathbf{s}_{t} \sim \mathcal{D}} \left[D_{\mathrm{KL}} \left(\pi_{\theta} \left(\cdot \mid \mathbf{s}_{t} \right) \| \frac{\exp(Q_{\phi}(\mathbf{s}_{t}, \cdot))}{Z_{\theta}(\mathbf{s}_{t})} \right) \right].$$

$$J_{\pi}(\theta) = \mathbb{E}_{\mathbf{s}_{t} \sim \mathcal{D}, \epsilon_{t} \sim \mathcal{N}} \left[\log \pi_{\theta} \left(f_{\theta} \left(\epsilon_{t}; \mathbf{s}_{t} \right) \mid \mathbf{s}_{t} \right) - Q_{\phi} \left(\mathbf{s}_{t}, f_{\theta} \left(\epsilon_{t}; \mathbf{s}_{t} \right) \right) \right],$$

$$\hat{\nabla}_{\theta} J_{\pi}(\theta) = \nabla_{\theta} \log \pi_{\theta} \left(\mathbf{a}_{t} \mid \mathbf{s}_{t} \right) + \left(\nabla_{\mathbf{a}_{t}} \log \pi_{\theta} \left(\mathbf{a}_{t} \mid \mathbf{s}_{t} \right) - \nabla_{\mathbf{a}_{t}} Q \left(\mathbf{s}_{t}, \mathbf{a}_{t} \right) \right) \nabla_{\theta} f_{\theta} \left(\epsilon_{t}; \mathbf{s}_{t} \right)$$

Replacing Prior $p(a_t|s_t)$

• Decision Choice 1: As we were not interested in action prior (online), we can take it as simplest:

$$p(a_t) = p(a_t|s_t) = \frac{1}{|\mathcal{A}_t|}$$

• What if we assume the prior as $p(a_t|s_t) = \pi_{old}(a_t|s_t)$ or $\pi_{behavior}(a_t|s_t)$?

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{T} r(s_t, a_t) + \alpha \mathcal{H} \left(\pi \left(\cdot \mid s_t \right) \right) \right]$$

Max Entropy objective leads to KL regularized objective. TRPO/PPO/AWR....

$$J(\pi) = \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{T} r(s_t, a_t) \right] - \alpha \mathbb{E}_{s \sim \rho_{\pi}} \left[D_{KL} \left(\pi(\cdot \mid s) || \pi_{\text{old}} \left(\cdot \mid s \right) \right) \right]$$

Trajectory planning. Diffuser

Modeling Entire trajectory distribution with offline dataset.

$$p(\tau|\mathcal{O}_{1:T}) = p(\tau) \exp\left(\sum_{t=1}^{T} r(s_t, a_t)\right)$$

Future Plan of Research

- 1. Constrained black-box optimization with posterior inference. (Taeyoung Yun, Kyuil Sim) ~NIPS(2025)
- 2. Probabilistic Inference for sequential decision making. (Diffusion meets control/Control meets diffusion)

EOD

