

# Paper summary: Stochastic latent actor-critic: deep reinforcement learning with a latent variable model

May 22, 2022

- 1 Idea in few sentences
- 2 Explanation of the central concept
- 3 Methodology
- 4 Initial rambly notes

## 4.1 Abstract

Stochastic latent actor-critic (SLAC): unifies stochastic sequential models and RL into a single method by learning a compact latent representation and then performing RL in the model’s learned latent space. In continuous actions spaces, on images, simultaneous training of representations and policy, latent representation and latent-space dynamics are learned jointly.

## 4.2 Introduction

A predictive model is trained. It combines learning stochastic sequential models and RL into a single model, performing RL in the model’s learned latent space. By formalizing the control problem as an inference problem within a POMDP, it can be shown that variational inference leads to the objective of the SLAC algorithm. Unlike model-based approaches which compound model error when planning, SLAC performs infinite horizon optimization.

### 4.2.1 Preliminary: maximum entropy RL in fully observable MDPs

MDP, states  $\mathbf{s}_t \in \mathcal{S}$ , actions  $\mathbf{a}_t \in \mathcal{A}$ , rewards  $r_t$ , initial state distribution  $p(\mathbf{s}_1)$ , stochastic transition distribution  $p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$ . Standard RL aims to learn parameters  $\phi$  under policy  $\theta_\phi(\mathbf{a}_t|\mathbf{s}_t)$  s.t. the expected sum of rewards is maximized

under the induced trajectory distribution  $\rho_\pi$ . This can be modified to include an entropy term, such that the policy also maximizes the expected entropy  $\mathcal{H}(\pi_\phi(\cdot|\mathbf{s}_t))$ . This is built on here. The resulting objective in maximum entropy RL is:

$$\sum_{t=1}^T \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} [r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi_\phi(\cdot|\mathbf{s}_t))] \quad (1)$$

where  $\alpha$  is the temperature parameter for balancing between reward and entropy contributions to the objective. SAC uses this to derive soft policy iteration. SAC has the critic  $Q_\theta$  and actor  $\pi_\phi$ . The soft Q-function parameters  $\theta$  are optimized to minimize the soft Bellman residual:

$$J_Q(\theta) = \frac{1}{2} (Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - (r_t + \gamma \mathbb{E}_{\mathbf{a}_{t+1} \sim \pi_\phi} [Q_{\bar{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \alpha \log \pi_\phi(\mathbf{a}_{t+1}|\mathbf{s}_{t+1})]))^2 \quad (2)$$

where  $\bar{\theta}$  are delayed parameters. The policy parameters  $\phi$  are optimized to update the policy towards the exponential of the soft Q-function, resulting in policy loss:

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} [\alpha \log(\pi_\phi(\mathbf{a}_t|\mathbf{s}_t)) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t)] \quad (3)$$

#### 4.2.2 Preliminary: sequential latent variable models and amortized variational inference in POMDPs

Latent variable models with amortized variational inference are used. Image is  $\mathbf{x}$ , generated latent representation is  $\mathbf{z}$ . The model is learned by maximizing the probability of each observed  $\mathbf{x}$  from some training set under the entire generative process

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \quad (4)$$

This is intractable to compute due to the marginalization of the latent variables  $\mathbf{z}$ . In amortized variational inference, the evidence lower bound for log-likelihood is used:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q} [\log p(\mathbf{z}|\mathbf{x})] - D_{KL}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (5)$$

$\log p(\mathbf{x})$  is maximized by learning an encoder  $q(\mathbf{z}|\mathbf{x})$  and a decoder  $p(\mathbf{x}|\mathbf{z})$  and directly performing gradient descent of the right hand side of the equation. Then the distributions of interest are the prior  $p(\mathbf{z})$ , the observation model  $p(\mathbf{x}|\mathbf{z})$  and the variational approximate posterior  $q(\mathbf{z}|\mathbf{x})$ .

To extend such models to sequential decision making, actions and temporal structure must be incorporated to the latent state. Since we're now in the POMDP setting, we don't have states, only observations  $\mathbf{x}_t \in \mathcal{X}$  and latent states  $\mathbf{z}_t \in \mathcal{Z}$ . The transition distributions are now  $p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{a}_t)$  and the observation model is  $p(\mathbf{x}_t|\mathbf{z}_t)$ . As VAEs, a generative model can be learned by maximizing the log-likelihood. Importantly,  $\mathbf{x}_t$  alone can not give  $\mathbf{z}_t$  — for that previous observations also must be taken into account. Hence the need for sequential latent variable models. The distributions of interest are  $p(\mathbf{z}_1)$  and  $p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{a}_t)$ , the observation model  $p(\mathbf{x}_t|\mathbf{z}_t)$  and the approximate variational

posteriors  $q(\mathbf{z}_1|\mathbf{x}_1)$  and  $q(\mathbf{z}_{t+1}|\mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t)$ . With these, the log-likelihood of the observations can be bounded:

$$\log p(\mathbf{x}_{1:\tau+1}|\mathbf{a}_{1:\tau}) \geq \mathbb{E}_{\mathbf{z}_{1:\tau+1} \sim q} \left[ \sum_{t=0}^{\tau} \log p(\mathbf{x}_{t+1}|\mathbf{z}_{t+1}) - D_{KL}(q(\mathbf{z}_{t+1}|\mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t) || p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{a}_t)) \right] \quad (6)$$

Where a few obvious conditionals have been dropped. Prior work has dealt with these non-Markovian state transitions with recurrent networks or probabilistic state-space models. Here, a stochastic latent variable model is trained.

### 4.3 Joint modelling and control as inference

In the MDP setting, the control problem can be embedded into the MDP graphical model by introducing a binary random variable  $\mathcal{O}_t$  which indicates if the time step  $t$  is optimal. If its distribution is chosen to be  $p(\mathcal{O}_t = 1|\mathbf{s}_t, \mathbf{a}_t) = \exp(r(\mathbf{s}_t, \mathbf{a}_t))$ , then the maximization of  $p(\mathcal{O}_{1:T})$  via approximate inference in that model yields the optimal policy for the maximum entropy objective. Now this is extended to the POMDP setting. Analogously, we get  $p(\mathcal{O}_t = 1|\mathbf{z}_t, \mathbf{a}_t) = \exp(r(\mathbf{z}_t, \mathbf{a}_t))$ . But now not only the likelihood of optimality variables is maximized, but also the maximum entropy policies, thus maximizing the marginal likelihood  $p(\mathbf{x}_{1:\tau+1}, \mathcal{O}_{\tau+1:T}|\mathbf{a}_{1:\tau})$ , in order to learn the model and the policy at the same time. This objective represents both the likelihood of the observed data from the past  $\tau + 1$  steps, as well as the optimality of the agent's action for future steps, hence combining both representation learning and control into a single graphical model.

The variational distribution is factorized into a product of *recognition* terms  $q(\mathbf{z}_{t+1}|\mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t)$ , *dynamics terms*  $p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{a}_t)$  and *policy* terms  $\pi(\mathbf{a}_t|\mathbf{x}_{1:t}, \mathbf{a}_{1:t-1})$ :

$$q(\mathbf{z}_{1:T}, \mathbf{a}_{\tau+1:T}|\mathbf{x}_{1:\tau+1}, \mathbf{a}_{1:\tau}) = \prod_{t=0}^{\tau} q(\mathbf{z}_{t+1}|\mathbf{x}_{t+1}, \mathbf{z}_t, \mathbf{a}_t) \prod_{t=\tau+1}^{T-1} p(\mathbf{z}_{t+1}|\mathbf{z}_t, \mathbf{a}_t) \prod_{t=\tau+1}^T \pi(\mathbf{a}_t|\mathbf{x}_{1:t}, \mathbf{a}_{1:t-1}) \quad (7)$$

The dynamics for future timesteps is used to prevent the agent from taking optimistic steps. This posterior can then be used to obtain the evidence lower bound (ELBO) of the likelihood. It is in turn separated into training the latent variable model and maximizing the entropy of RL. I'm not going into this as I can't understand it without much more effort. The latent model is a specific VAE and I don't exactly understand the vanilla VAE anyway. Also I never derived SAC and I don't know anything about messaging passing in statistics so I certainly can't follow how the maximum entropy actor-critic is employed here. Let's call it good because unlike other more complex generative models it actually integrates them into the POMDP (which is the case, but I have to take it on belief here).

**4.4 Method**

**4.5 Other stuff**