

# Paper summary: Improving sample efficiency in model-free reinforcement learning from images

May 29, 2022

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

## **4.1 Abstract**

Fitting a high-capacity encoder to extract features (state information) from images with only the reward signal leads to poor performance. One option is to incorporate reconstruction loss into an off-policy algorithm, but that often leads to training instability. Investigation into why shows variational autoencoders to be a problem.

## **4.2 Introduction**

Some solutions to low sample efficiency are:

1. use an off-policy algorithm
2. add an auxiliary task with an unsupervised objective

The simplest auxiliary task is an autoencoder with a pixel reconstruction objective. Prior works use a two-step training procedure, but this often leads to lower final performance.

We confirm that a pixel reconstruction loss is vital for learning a good representation, specifically when trained jointly, but requires careful design choices to succeed.

There are 3 contributions:

1. methodical study of the issues involved with combining autoencoders with model-free RL in the off-policy setting, resulting in SAC+AE
2. demonstrating SAC+AE does its thing robustly
3. open-source code of SAC+AE

### 4.3 Related work

The 2010 AE paper re-encodes after every AE update, which is unfeasible for large problems. In Finn et al. in Learning visual feature spaces for robotic manipulation with deep spatial autoencoders, authors pretrain and the linear policy is trained separately. This does not translate to end-to-end methods which are to be developed here. Other stuff was unstable and hindered policy learning performance. Model-based stuff is cool and is efficient, but it is super brittle and sensitive to hyperparameters due to multiple different auxiliary losses, ex. dynamics loss, reward loss, decoder loss, in addition to policy and/or value optimizations.

### 4.4 Background

#### 4.4.1 SAC

Maximum entropy objective:

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

This is used to derive soft policy iteration. Soft Bellman residual:

$$J(Q) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1}) \sim \mathcal{D}} \left[ \left( Q(\mathbf{s}_t, \mathbf{a}_t) - r_t - \gamma \bar{V}(\mathbf{s}_{t+1}) \right)^2 \right] \quad (2)$$

where the target value function  $\bar{V}$  is approximate via Monte-Carlo estimate of

$$\bar{V}(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [\bar{Q}(\mathbf{s}_t, \mathbf{a}_t) - \alpha \log \pi(\mathbf{a}_t | \mathbf{s}_t)] \quad (3)$$

The policy improvement step then attempts to project a parametric policy  $\pi(\mathbf{a}_t | \mathbf{s}_t)$  by minimizing KL divergence between the policy and a Boltzmann distribution induced by the Q-function using the objective:

$$J(\pi) = \sim \sim \mathcal{D} [D_{KL}(\pi(\cdot | \mathbf{s}_t) || Q(\mathbf{s}_t, \cdot))] \quad (4)$$

where  $Q(\mathbf{s}_t, \cdot) \propto \exp \left\{ \frac{1}{\alpha} Q(\mathbf{s}_t, \cdot) \right\}$

#### 4.4.2 Image-based observations and autoencoders

AE is represented as a convolutional encoder  $g_\phi$  that maps an image observation  $\mathbf{o}_t$  to a low-dimensional latent vector  $\mathbf{z}_t$ , and a deconvolutional decoder  $f_\phi : \mathcal{Z} \rightarrow \mathcal{O}$ . Both the encoder and decoder are trained simultaneously by maximizing the expected log-likelihood

$$J(AE) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\log p_\theta(\mathbf{o}_t | \mathbf{z}_t)] \quad (5)$$

where  $\mathbf{z}_t = g_\phi(\mathbf{o}_t)$ .

In the  $\beta$ -VAE case, the objective is:

$$J(VAE) = \mathbb{E}_{\mathbf{o}_t \sim \mathcal{D}} [\mathbb{E}_{\mathbf{z}_t \sim q_\phi(\mathbf{z}_t | \mathbf{o}_t)} [\log p_\theta(\mathbf{o}_t | \mathbf{z}_t)]] - \beta D_{KL}(q_\phi(\mathbf{z}_t | \mathbf{o}_t) || p(\mathbf{z}_t)) \quad (6)$$

where the variational distribution is parametrized as  $q_\phi(\mathbf{z}_t | \mathbf{o}_t) = \mathcal{N}(\mathbf{z}_t | \mu_\phi(\mathbf{o}_t), \sigma_\phi^2(\mathbf{o}_t))$ . The latent vector  $\mathbf{z}_t$  is used by an RL algorithm instead of the unavailable true state  $\mathbf{s}_t$ .

### 4.5 Representation learning with image reconstruction

For a model-free RL algorithm, learning from pixels yield much worse results than learning from state. Prior works have shown that using auxiliary supervision to learn state representations helps. The focus of this work is to examine the use of image reconstruction loss as the auxiliary loss. Task-dependent auxiliary loss and world models are avoided in this work.

The authors tried to use a  $\beta$ -VAE, but only on current frames, instead of a sequence of frames. They tried alternating the training, and observed a positive correlation between performance and alternation frequency, but the final result didn't close the performance gap. They tried updating the  $\beta$ -VAE encoder with actor-critic gradients, but this led to severe instability in training. They conclude that they resulted from the stochastic nature of  $\beta$ -VAEs and the non-stationary gradient from the actor.

#### 4.5.1 Alternating representation learning with a $\beta$ -VAE

First thing is to confirm that alternative training between the AE and RL algorithm.

### 4.6 Method

#### 4.7 Other stuff