

---

# State Abstraction as Compression in Apprenticeship Learning

## Supplementary Material

---

**David Abel<sup>1</sup>, Dilip Arumugam<sup>2</sup>, Kavosh Asadi<sup>1</sup>**  
**Yuu Jinnai<sup>1</sup>, Michael L. Littman<sup>1</sup>, Lawson L.S. Wong<sup>3</sup>**

1: Department of Computer Science, Brown University.  
 2: Department of Computer Science, Stanford University.  
 3: College of Computer and Information Science, Northeastern University.

We here include proofs of the introduced theorems and lemmas and provide further discussion about extensions to more general settings.

### Proofs

We begin with the two central Lemmas, which serve as the main deductive steps for the proof of Theorem 2.

.....

**Lemma 1.** *Consider a discrete random variable  $X$ , with alphabet  $\mathcal{X}$  and some pmf  $p(x)$ . For a given threshold  $\delta_{min} \in (0, 1)$ , the pmf-used alphabet size of the alphabet is bounded:*

$$|\mathcal{X}|_{p(x)}^{\delta_{min}} \leq \frac{H(X)}{\delta_{min} \log\left(\frac{1}{\delta_{min}}\right)}. \quad (18)$$

*Proof.* Recall the entropy of a random variable  $H(X)$ :

$$H(X) := - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x). \quad (19)$$

Then, for a given maximum entropy  $H(X) \leq N$ , we seek to understand the largest possible  $|\mathcal{X}|_{p(x)}^{\delta_{min}}$ . That is:

$$\max_{p(x): H(X) \leq N} |\mathcal{X}|_{p(x)}^{\delta_{min}}. \quad (20)$$

Note that this is maximized at the uniform distribution, where each element has  $\delta_{min}$  probability, across the largest alphabet such that  $H(X) = N$ :

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x) \quad (21)$$

$$= - \sum_{x \in \mathcal{X}} \delta_{min} \log_2 \delta_{min} \quad (22)$$

$$= -|\mathcal{X}| \delta_{min} \log_2 \delta_{min} \quad (23)$$

$$= |\mathcal{X}| \delta_{min} \log_2 \frac{1}{\delta_{min}}. \quad (24)$$

Therefore, for a given pmf  $p(x)$  with entropy  $H(X)$ , and a minimum threshold of probability, the minimum size of the alphabet  $\mathcal{X}$  is upper bounded:

$$|\mathcal{X}| \leq \frac{H(X)}{\delta_{min} \log_2 \frac{1}{\delta_{min}}}. \quad \square$$

**Lemma 2.** Consider two stochastic policies,  $\pi_1$  and  $\pi_2$  on state space  $\mathcal{S}$ , and a fixed probability distribution over  $\mathcal{S}$ ,  $p(s)$ . If, for some  $k \in \mathbb{R}_{\geq 0}$ :

$$\mathbb{E}_{p(s)} [D_{KL}(\pi_1(a | s) || \pi_2(a | s))] \leq k, \quad (25)$$

then:

$$\mathbb{E}_{p(s)} [V^{\pi_1}(s) - V^{\pi_2}(s)] \leq \sqrt{2k} \text{VMAX}. \quad (26)$$

*Proof.* Recall the total variation distance (TVD) between our two policies for a given state  $s$  is defined as:

$$TV(\pi_E, \pi_\phi) := \sup_{a \in \mathcal{A}} |\pi_E(a | s) - \pi_\phi(a | s)|. \quad (27)$$

Furthermore, recall that TVD relates to the  $L_1$  norm and the KL divergence:

$$TV(\pi_E, \pi_\phi) = \frac{1}{2} \sum_{a \in \mathcal{A}} |\pi_E(a) - \pi_\phi(a)| \leq \sqrt{\frac{1}{2} D_{KL}(\pi_E || \pi_\phi)}, \quad (28)$$

where the inequality in Equation 28 is formally known as Pinsker's inequality. With this inequality in place, we expand the expectation in the value bound:

$$\begin{aligned} & \mathbb{E}_{p(s)} [V^{\pi_E}(s) - V^{\pi_\phi}(s)] \quad (29) \\ & \leq \sum_s p(s) \left( \sum_a |\pi_E(a | s) - \pi_\phi(a | s)| (\mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{T}(s, a, s') |V^{\pi_E}(s') - V^{\pi_\phi}(s')|) \right) \\ & = \sum_s \sum_a p(s) |\pi_E(a | s) - \pi_\phi(a | s)| \left( \mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{T}(s, a, s') |V^{\pi_E}(s') - V^{\pi_\phi}(s')| \right) \end{aligned}$$

Then, applying the upper bound on the possible value  $\text{VMAX} = \text{RMAX}/(1 - \gamma)$  to Equation 29:

$$\mathbb{E}_{p(s)} [V^{\pi_E}(s) - V^{\pi_\phi}(s)] \leq \text{VMAX} \mathbb{E}_{p(s)} [|\pi_E(a | s) - \pi_\phi(a | s)|]. \quad (30)$$

Then, by Pinsker's inequality, we conclude:

$$\begin{aligned} \mathbb{E}_{p(s)} [V^{\pi_E}(s) - V^{\pi_\phi}(s)] & \leq 2\text{VMAX} \mathbb{E}_{p(s)} \left[ \sqrt{\frac{1}{2} D_{KL}(\pi_E(a | s) || \pi_\phi(a | s))} \right] \quad (31) \\ & \leq \sqrt{2k} \text{VMAX} \quad \square \end{aligned}$$

With these two Lemmas in place, we now prove our main Theorem.

**Theorem 2.** A function  $f$  of the DIB objective  $\hat{\mathcal{J}}_{DIB}$  is an upper bound for the CVA Objective,  $\mathcal{J}$ , where state space size is treated as  $|\mathcal{S}_\phi|_{\rho_\phi(s)}^{\delta_{min}}$ :

$$\forall_\phi : \mathcal{J}[\phi] \leq f \left( \hat{\mathcal{J}}_{DIB}[\phi] \right). \quad (32)$$

*Proof.* With Lemma 1 and Lemma 2, the proof is straightforward. Consider the  $\phi$  that minimizes  $\hat{\mathcal{J}}_{\text{DIB}}$ , yielding the value of at most  $N + \beta k$ , where:

$$N := H(S_\phi) \tag{33}$$

$$k := \mathbb{E}_{s \sim \rho_E} [D_{KL}(\pi_E(a | s) || \pi_\phi(a | \phi(s)))]. \tag{34}$$

Then, by Lemma 1, we know:

$$|\mathcal{S}_\phi|_{\rho_\phi(s)}^{\delta_{min}} \leq \frac{N}{\delta_{min} \log_2 \frac{1}{\delta_{min}}}. \tag{35}$$

By Lemma 2, we know:

$$\mathbb{E}_{\rho(s)} [V^{\pi_E}(s) - V^{\pi_\phi}(s)] \leq \sqrt{2k} \text{VMAX}. \tag{36}$$

Therefore, since both quantities are non-negative, we conclude:

$$\mathcal{J}[\phi] = |\mathcal{S}_\phi|_{\rho_\phi(s)}^{\delta_{min}} + \beta \mathbb{E}_{\rho(s)} [V^{\pi_E}(s) - V^{\pi_\phi}(s)] \tag{37}$$

$$\leq \frac{N}{\delta_{min} \log_2 \frac{1}{\delta_{min}}} + \beta \sqrt{2k} \text{VMAX}. \tag{38}$$

Thus, we can upper bound the quantities in  $\mathcal{J}$  as a function of the quantities in  $\hat{\mathcal{J}}_{\text{DIB}}$ . □

.....

## Extensions

The results we present in this work offer a first step toward understanding state abstraction as compression in sequential decision making. Looking forward, we hope to extend this work into a general abstraction-as-compression learning paradigm for RL. To this end, we introduce two concrete avenues for extending the work to more general settings.

### Agent Has Control

Relaxing the assumption that  $\pi_E$  controls the MDP is essential for extending this work to traditional RL. We here propose a path toward removing this restriction by focusing on an intermediate goal: define an algorithm with the same properties as DIBS, but with the learning agent’s *non-stationary* policy controlling the underlying MDP instead of  $\pi_E$ . Ultimately, we seek an algorithm that, after  $T < \infty$  iterations, can produce an abstraction–policy pair such that, for some state distribution  $p(s)$ :

$$\mathbb{E}_{p(s)} [V^d(s) - V^{\pi_\phi^T}(s)] \leq f(\beta, T). \tag{39}$$

The most challenging aspect of this setup is that the source distribution is no longer fixed, since the agent’s policy will change over time as the agent learns and updates both  $\phi$  and  $\pi_\phi$ . To this end, we prove the following result that suggests a route to defining a convergent algorithm for the case where the agent controls the MDP.

**Lemma 3.** *Given two policies  $\pi_1$  and  $\pi_2$ , if:*

$$\sup_s \sum_a |\pi_1(a | s) - \pi_2(a | s)| \leq \Delta, \tag{40}$$

*then:*

$$\sum_s |\rho_{\pi_1, s_0}(s) - \rho_{\pi_2, s_0}(s)| \leq \frac{\Delta \gamma}{1 - \gamma}, \tag{41}$$

*where  $\rho_{\pi, s_0}$  denotes the stationary distribution over states under  $\pi$ , starting in state  $s_0$ .*

*Proof.* We first bound the difference between the two state distributions after  $t$  steps:

$$\begin{aligned}
& \sum_{s'} |\rho_{\pi_1, s_0}^t(s') - \rho_{\pi_2, s_0}^t(s')| \\
= & \sum_{s'} |p(s_t = s' | s_0, \pi_1) - p(s_t = s' | s_0, \pi_2)| \\
= & \sum_{s'} \left| \sum_s p(s_{t-1} = s | s_0, \pi_1) \sum_a \pi_1(a | s) p(s' | s, a) \right. \\
& \left. - \sum_s p(s_{t-1} = s | s_0, \pi_2) \sum_a \pi_2(a | s) p(s' | s, a) \right| \\
\leq & \sum_{s'} \left| \sum_s p(s_{t-1} = s | s_0, \pi_1) \sum_a \left( \pi_1(a | s) - \pi_2(a | s) \right) p(s' | s, a) \right| \\
& + \sum_{s'} \left| \sum_s \left( p(s_{t-1} = s | s_0, \pi_1) - p(s_{t-1} = s | s_0, \pi_2) \right) \sum_a \pi_2(a | s) p(s' | s, a) \right| \\
\leq & \sum_s p(s_{t-1} = s | s_0, \pi_1) \sum_a \left| \pi_1(a | s) - \pi_2(a | s) \right| \sum_{s'} p(s' | s, a) \\
& + \sum_s \left| p(s_{t-1} = s | s_0, \pi_1) - p(s_{t-1} = s | s_0, \pi_2) \right| \sum_a \pi_2(a | s) \sum_{s'} p(s' | s, a) \\
\leq & \Delta + \sum_s \left| p(s_{t-1} = s | s_0, \pi_1) - p(s_{t-1} = s | s_0, \pi_2) \right| \\
= & \Delta + \sum_{s'} |\rho_{\pi_1, s_0}^{t-1}(s') - \rho_{\pi_2, s_0}^{t-1}(s')|
\end{aligned}$$

From the above bound, and using induction, we have:

$$\begin{aligned}
& \sum_{s'} |\rho_{\pi_1, s_0}^t(s') - \rho_{\pi_2, s_0}^t(s')| \leq t\Delta \\
& \sum_s |\rho_{\pi_1, s_0}(s) - \rho_{\pi_2, s_0}(s)| \\
= & \sum_s \left| (1 - \gamma) \sum_t \gamma^t \rho_{\pi_1, s_0}^t(s) - (1 - \gamma) \sum_t \gamma^t \rho_{\pi_2, s_0}^t(s) \right| \\
\leq & (1 - \gamma) \sum_t \gamma^t \sum_s |\rho_{\pi_1, s_0}^t(s) - \rho_{\pi_2, s_0}^t(s)| \\
\leq & (1 - \gamma) \sum_t \gamma^t t \Delta = (1 - \gamma) \frac{\gamma \Delta}{(1 - \gamma)^2} = \frac{\gamma \Delta}{1 - \gamma} \quad \square
\end{aligned}$$

.....

**Corollary 1.** As a simple corollary of Proposition 1 we get:

$$\begin{aligned}
|V^{\pi_1}(s) - V^{\pi_2}(s)| &= \left| \sum_t \sum_s \rho_{\pi_1, s_0}^t(s) R(s) - \sum_t \sum_s \rho_{\pi_2, s_0}^t(s) R(s) \right| \\
&\leq \text{RMAX} \sum_t (1 - \gamma) \sum_s |\rho_{\pi_1, s_0}^t(s) - \rho_{\pi_2, s_0}^t(s)| \\
&\leq \text{RMAX} (1 - \gamma) \sum_t \Delta \gamma^t t \\
&= \text{RMAX} \frac{\gamma \Delta}{1 - \gamma} = \Delta \gamma \text{VMAX}
\end{aligned}$$

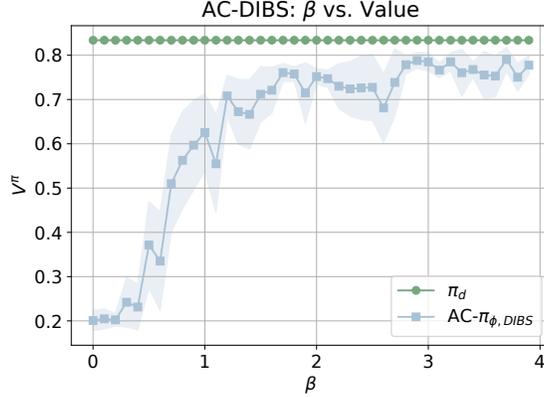


Figure 1: The value of the abstract policy found by Agent Controlled DIBS for values of  $\beta$  between 0 and 4. The line indicates the average over 100 trials, with 95% confidence intervals.

The above Lemma suggests that two policies that deviate by a bounded amount are guaranteed to share similar stationary distributions. So, with the right intermediate updates that encourage the abstract policy to mimic  $\pi_E$  more faithfully, we can construct a convergent algorithm for the agent-in-control setting.

We offer an initial variant of this algorithm, which we called agent-controlled DIBS (AC-DIBS). We conduct a Four Rooms experiment similar to that of the previous section. Here, we now run the entire process of DIBS for  $T$  rounds, each time to convergence, but letting the agent’s *initial* policy for that round define the stationary state distribution. Results are presented in Figure 1. Surprisingly, we find AC-DIBS *always* converges quickly to a point where  $\rho_\phi^T$  is sufficiently close to  $\rho_E$ , when  $\beta > 1$ . This finding supports our supposition, along with Lemma 3, there is a route to defining a convergent form of DIBS when the agent’s policy controls the underlying MDP.

## Multiple MDPs

Second, we use our framework to compute a single abstraction sufficient for representing the demonstrator policy across different, but potentially related tasks. We suppose we are given a set of MDPs  $\mathcal{M}$ , each sharing a state and action space, but are allowed to vary in  $\mathcal{T}$ ,  $\mathcal{R}$ , and  $\gamma$ . We conduct an experiment in which  $|\mathcal{M}| = 4$ , each with a goal in one of the four corners of the world. We run DIBS for each MDP in  $\mathcal{M}$ , for a fixed  $\beta$ , and form a master abstraction  $\phi_{\mathcal{M}}$  by taking the intersection across each computed state abstraction. That is, for any state pair  $(s_1, s_2)$ , for  $\phi_i$  computed by DIBS on each MDP, we define  $\phi_{\mathcal{M}}$ :

$$\phi_{\mathcal{M}}(s_1) = \phi_{\mathcal{M}}(s_2) \equiv \bigwedge_{i=1}^{|\mathcal{M}|} \{\phi_i(s_1) = \phi_i(s_2)\}. \quad (42)$$

Figure 2 shows  $\phi_{\mathcal{M}}$  for different values of  $\beta$ . All cells with the same color are grouped into the same state, except for white: all white states are each treated as their true ground state. Note that the abstraction becomes far more detailed as  $\beta$  increases: when  $\beta$  is close to 0, the algorithm prioritizes *compression*, as is reflected by Figure 2a, which only has a single state. Conversely, as  $\beta$  increases, we find the algorithm adds more distinctions between states, only grouping those that are close to one another or the near the same wall. When  $\beta = 1$ , we find that the abstraction groups the top and bottom hallways together, and the left and right hallways together. Similarly, it groups large regions of contiguous states together, such as the central group of states in the bottom right room, and the right most states in the top right room. When  $\beta = 10$ , we find less compression, but still find contiguous regions that have some structural similarities. For instance, the dark blue regions extending above and below the left and right doorways are grouped together, as across all four goals it is important to move up and down between the rooms. The orange groups in the top left and bottom right rooms suggest a similar structure for moving left and right. Critically, none of the abstractions are perfect: we do not know what constitutes the optimal abstraction in the case, nor how well our proposed algorithm approximates this optimal. In future work we hope to bring clarifying answers to these questions.

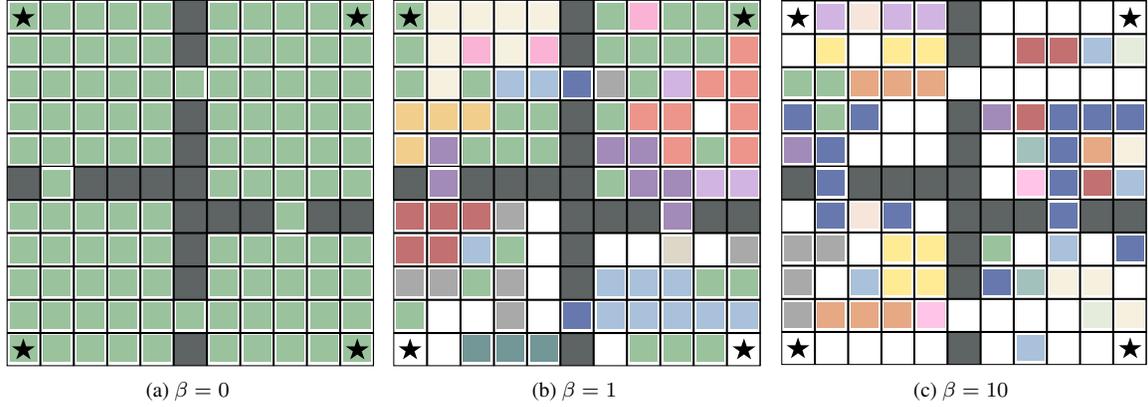


Figure 2: State abstractions computed by DIBS for a collection of MDPs for different values of  $\beta$ .

.....

Lastly, we briefly discuss the variational upper bound introduced as part of our VAE extension. In Equation 14, we reference a known result relating the KL-divergence to mutual information in the context of variational autoencoders (VAEs) Kingma and Welling (2013). Recall that a VAE is concerned with learning a probabilistic encoder-decoder pair,  $q_\psi(z|x)$  and  $p_\theta(x|z)$  respectively, for compressing raw data into a latent representation and then reconstructing the original input. During the course of training a VAE via the evidence lower bound objective (ELBO), the KL-divergence between  $q_\psi(z|x)$  and some prior distribution over latent codes,  $p(z)$ , is minimized in expectation over the data distribution,  $p(x)$ ; that is,  $\mathbb{E}_{p(x)}[D_{KL}(q_\psi(z|x) || p(z))]$  is minimized.

To define our variational upper bound to  $\hat{J}$  (the stochastic IB objective for state abstraction), we leverage a known result Makhzani and Frey (2017); Kim and Mnih (2018); Dupont (2018) whose proof we replicate here for completeness with  $q(z, x) = p(x)q_\psi(z|x)$  and  $q(z) = \mathbb{E}_{p(x)}[q_\psi(z|x)]$ :

$$\begin{aligned}
 \mathbb{E}_{p(x)}[D_{KL}(q_\psi(z|x) || p(z))] &= \mathbb{E}_{p(x)} \left[ \mathbb{E}_{q_\psi(z|x)} \left[ \log \frac{q_\psi(z|x)}{p(z)} \right] \right] \\
 &= \mathbb{E}_{q(z,x)} \left[ \log \left( \frac{q_\psi(z|x) q(z)}{p(z) q(z)} \right) \right] \\
 &= \mathbb{E}_{q(z,x)} \left[ \log \frac{q_\psi(z|x)}{q(z)} \right] + \mathbb{E}_{q(z,x)} \left[ \log \frac{q(z)}{p(z)} \right] \\
 &= \mathbb{E}_{q(z,x)} \left[ \log \frac{q(z,x)}{q(z)p(x)} \right] + \mathbb{E}_{q(z)} \left[ \log \frac{q(z)}{p(z)} \right] \\
 &= I(X; Z) + D_{KL}(q(z) || p(z)) \\
 &\geq I(X; Z)
 \end{aligned}$$

.....

## References

- Dupont, E. 2018. Joint-vae: Learning disentangled joint continuous and discrete representations. *arXiv preprint arXiv:1804.00104*.
- Kim, H., and Mnih, A. 2018. Disentangling by Factorising. *arXiv preprint arXiv:1802.05983*.
- Kingma, D. P., and Welling, M. 2013. Auto-Encoding Variational Bayes. *Proceedings of the International Conference on Learning Representations*.

Makhzani, A., and Frey, B. J. 2017. Pixelgan autoencoders. In *Advances in Neural Information Processing Systems*, 1975–1985.