# Interpretable Option Discovery using Deep Q-Learning and Variational Autoencoders

Per-Arne Andersen[✉], Morten Goodwin, and Ole-Christoffer Granmo

Department of ICT, University of Agder, Grimstad, Norway
{per.andersen,morten.goodwin,ole.granmo}@uia.no

**Abstract.** Deep Reinforcement Learning (RL) is unquestionably a robust framework to train autonomous agents in a wide variety of disciplines. However, traditional deep and shallow model-free RL algorithms suffer from low sample efficiency and inadequate generalization for sparse state spaces. The options framework with temporal abstractions [18] is perhaps the most promising method to solve these problems, but it still has noticeable shortcomings. It only guarantees local convergence, and it is challenging to automate initiation and termination conditions, which in practice are commonly hand-crafted.
Our proposal, the Deep Variational Q-Network (DVQN), combines deep generative- and reinforcement learning. The algorithm finds good policies from a Gaussian distributed latent-space, which is especially useful for defining options. The DVQN algorithm uses MSE with KL-divergence as regularization, combined with traditional Q-Learning updates. The algorithm learns a latent-space that represents good policies with state clusters for options. We show that the DVQN algorithm is a promising approach for identifying initiation and termination conditions for option-based reinforcement learning. Experiments show that the DVQN algorithm, with automatic initiation and termination, has comparable performance to Rainbow and can maintain stability when trained for extended periods after convergence.

**Keywords:** Deep Reinforcement Learning · Clustering · Options · Hierarchical Reinforcement Learning · Latent-space representation

## 1 Introduction

The interest in deep Reinforcement Learning (RL) is rapidly growing due to significant progress in several RL problems [2]. Deep RL has shown excellent abilities in a wide variety of domains, such as video games, robotics and, natural language progressing [16,14,13]. Current trends in applied RL has been to treat neural networks as black-boxes without regard for the latent-space structure. While unorganized latent-vectors are acceptable for model-free RL, it is disadvantageous for schemes such as options-based RL. In an option-based RL, the policy splits into sub-policies that perform individual behaviors based on the current state of the agent. A sub-policy, or option, is selected with initialization criteria and ended with a termination signal. The current state-of-the-art

in option-based RL primarily uses hand-crafted options. Option-based RL algorithms work well for simple environments but have poor performance in more complicated tasks. There is, to the best of our knowledge, no literature that addresses good option selection for difficult control tasks. There are efforts for making automatic options selection [17], but no method achieves notable performance across various environments.

This paper proposes a novel deep learning architecture for Q-learning using variational autoencoders that learn to organize similar states in a vast latent-space. The algorithm derives good policies from a latent-space that feature interpretability and the ability to classify sub-spaces for automatic option generation. Furthermore, we can produce human-interpretable visual representations from latent-space that directly reflects the state-space structure. We call this architecture DVQN for deep Variational Q-Networks and study the learned latent-space on classic RL problems from the Open AI gym [4].

The paper is organized as follows. Section 3 introduces preliminary literature for the proposed algorithm. Section 4 presents the proposed algorithm architecture. Section 5 outlines the experiment setup and presents empirical evidence of the algorithm performance. Section 2 briefly surveys work that is similar to our contribution. Finally, Section 6 summarises the work of this paper and outlines a roadmap for future work.

## 2   Related Work

There are numerous attempts in the literature to improve interpretability with deep learning algorithms, but primarily in the supervised cases. [22] provides an in-depth survey of interpretability with Convolutional Neural Networks (CNNs). Our approach is similar to the work of [20], where the authors propose an architecture for visual perception of the DQN algorithm. The difference, however, is primarily our focus on the interpretability of the latent-space distribution via methods commonly found in variational autoencoders. There are similar efforts to combine Q-Learning with Variational Autoencoders, such as [19,11], and shows promising results theoretically but with limited focus on interpretability. [1] did notable work on interpretability among using a distance KL-distance for optimization but did not find convincing evidence for a deeper understanding of the model. The focus of our contribution deviates here and finds significant value in a shallow and organized latent-space.

**Options** The learned latent-space is valuable for the selection of options in hierarchical reinforcement learning (HRL). There is increasing engagement in HRL research because of several appealing benefits such as sample efficiency and model simplicity [3]. Despite its growing attention, there are few advancements within this field compared to model-free RL. The options framework [18] is perhaps the most promising approach for HRL in terms of intuitive and convergence guarantees. Specifically, the options framework defines semi-Markov decision processes (SMDP), which is an extension of the traditional MDP framework [21]. SMDP features temporal abstractions where multiple discrete time steps are

generalized to a single step. These abstract steps are what defines an option, where the option is a subset of the state-space. In the proposed algorithm, the structure of the latent-space forms such temporal abstractions for options to form.

## 3   Background

The algorithm is formalized under conventional Markov decision processes tuples $< S, A, P, R, \gamma >$ where $S$ is a (finite) set of all possible states, $A$ is a (finite) set of all possible actions, $P$ defines the probabilistic transition function $P(S_{t+1} = s'|s, a)$ where $s$ is the previous state, and $s'$ is the transition state. $R$ is the reward function $R(r_{t+1}|s, a)$. Finally, the $\gamma$ is a discount factor between $\gamma \in [0 \dots 1]$ that determines the importance of future states. Lower $\gamma$ values decrease future state importance while higher values increase.
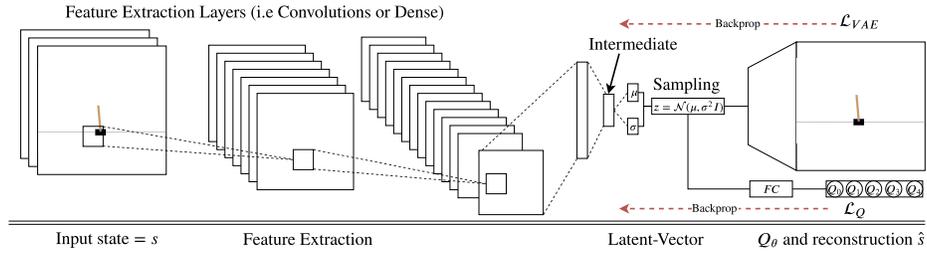
## 4   Deep Variational Q-Networks

Our contribution is a deep Q-learning algorithm that finds good policies in an organized latent space from variational autoencoders. [1] Empirically, the algorithm shows comparable performance to traditional model-free deep Q-Networks variants. We name our method the **D**eep **V**ariational **Q**-**N**etwork (DVQN) that combines two emerging algorithms, the variational autoencoder (VAE) [12] and deep Q-Networks (DQN) [14].

In traditional deep Q-Networks, the (latent-space) hidden layers are treated as a black-box. On the contrary, the objective of the variational autoencoder is to reconstruct the input and **organize** the latent-vector so that similar (data) states are adjacently modeled as a Gaussian distribution.

In DQN, the latent-space is sparse and is hard to interpret for humans and even option-based machines. By introducing a VAE mechanism into the algorithm, we expect far better interpretability for creating options in RL, which is the primary motivation for this contribution. Variational autoencoders are, in contrast to deep RL, involved with the organization of the latent-space representation, and commonly used to generate clusters of similar data with t-SNE or PCA [23]. The DVQN algorithm introduces three significant properties. First, the algorithm fits the data as a Gaussian distribution. This reduces the policy-space, which in practice reduces the probability of the policy drifting away from global minima. Second, the algorithm is generative and does not require exploration schemes such as $\epsilon$-greedy because it is done in re-parametrization during training. Third, the algorithm can learn the transition function and, if desirable, generate training data directly from the latent-space parameters, similar to the work of [8].

Figure 1 illustrates the architecture of the algorithm. The architecture follows general trends in similar RL literature but has notable contributions. First,

---

[1] The code will be published upon publication.

**Fig. 1.** The deep variational Q-Networks architecture.

features are extracted from the state-input, typically by using convolutions for raw images and fully-connected for vectorized input. The extracted features are forwarded to a fully connected intermediate layer of a user-specified size commonly between 50 to 512 neurons. The intermediate layer splits into two streams that represent the variance $\mu$ and standard deviation $\sigma$ and is used to sample the latent-vector using a Gaussian distribution through the re-parameterization. The latent-vector is forwarded to the decoder for state reconstruction and the Q-Learning stream for action-value (Q-value) optimization. The decoder and Q-Learning streams have the following loss functions:

$$\mathcal{L}_{VAE} = MSE(s, \hat{s}) + D_{KL}[q_\psi(z|s)\|p_\theta(z|s)] \tag{1}$$

$$\mathcal{L}_{DQN} = (r + \gamma Q(s', \arg\max_{a'} Q(s', a'; \theta_i); \theta_i) - Q(s, a; \theta_i))^2 \tag{2}$$

$$\mathcal{L}_{DVQN} = c_1 \mathbb{E}_{\sim q_\psi(z|s)}[\mathcal{L}_{VAE}] + c_2 \mathbb{E}_{s,a,s',D\sim r}[\mathcal{L}_{DQN}]. \tag{3}$$

The global loss function $\mathcal{L}_{DVQN}$ is composed of two local objectives: $\mathcal{L}_{DQN}$ and $\mathcal{L}_{VAE}$. In the VAE loss, the first term is the mean squared error between the input $s$ and its reconstruction $\hat{s}$. The second term is regularization using KL-distance to minimize the distance between the latent distribution and a Gaussian distribution. The DQN loss is a traditional deep Q-Learning update, as described in [14].

---

**Algorithm 1** DVQN: Minimal Implementation

---

1: Initialise $\Omega$
2: Initialise DVQN model $\pi$
3: Initialise replay-buffer $D_\pi$
4: **for** N episodes **do**
5:     $D_\pi \leftarrow$ Collect samples from $\Omega$ under the untrained policy $\pi$ via the generative policy sampling.
6:     Train model $\pi$ on a mini-batch from $D_\pi$ with objective from Equation 3

---

Algorithm 1 shows a general overview of the algorithm. First, the environment is initialized. Second, the DVQN model from Figure 1 is initialized with the

desired hyperparameters, and third, the replay-buffer is created. For a specified number of episodes, the algorithm samples actions from the generative policy for exploration and stores these as MDP tuples in the experience replay. After each episode, the algorithm samples mini-batches from the experience replay and performs parameter updates using stochastic gradient descent. The (loss function) optimization objective is described in equation 3. The process repeats until the algorithm converges.
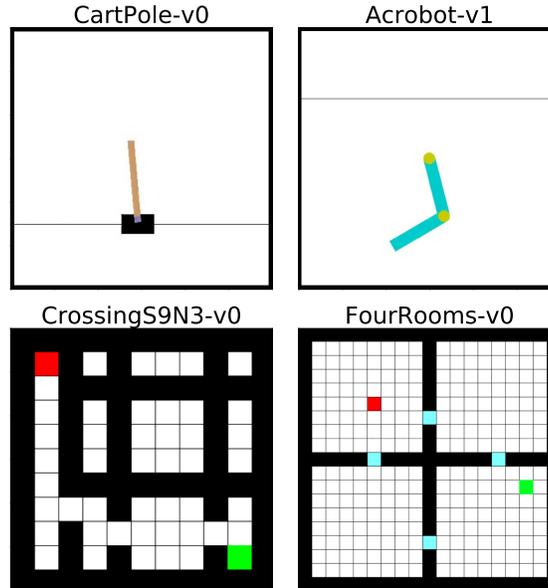
## 5    Experiments and Results

In this section, we conduct experiments against four traditional environments to demonstrate the effectiveness of the DVQN algorithm. We show that the algorithm can organize the latent-space by state similarity while maintaining comparable performance to model-free deep Q-learning algorithms.

### 5.1    Experiment test-bed

We evaluate the DVQN in the following environments; CartPole-v0, Acrobot-v1, CrossingS9N3-v0, and FourRooms-v0, shown in Figure 2. These environments are trivial to solve using model-free reinforcement learning and hence, excellent for visualizing the learned latent-space. The FourRooms-v0 environment is especially suited for option-based RL and can solve the problem in a fraction of time steps compared to model-free RL. Although the DVQN algorithm **does not quantify options for analysis** (see Section 6), the primary goal is to organize the latent-space so that it is possible to extract meaningful and interpretable options automatically. The aim is to have **comparable performance** to vanilla deep Q-Network variants found in the literature [14,5,10]. DVQN benchmarks against vanilla DQN, Double DQN, and Rainbow.

FourRooms-v0 and CrossingS9N3-v0 are a grid-world environment where the objective is to reach the terminal-state cell (In the lower right of the image in both environments). In FourRooms-v0, the agent has to enter several doors and only complete a part of the goal for each door it enters. FourRooms-v0 is an ideal environment for option-based reinforcement-learning because each door is considered a sub-goal. While the environment is solvable by many deep reinforcement learning algorithms, option-based RL is more efficient. The agent receives small negative rewards for moving and positive rewards for entering the goal-state (global) or the doors (local). The Crossing is a simpler environment where the agent has to learn the shortest path to the goal state. In both grid-environments, the agent can move in any direction, one cell per time step.

To further show that the algorithm works in simple control tasks, we perform experiments in CartPole-v0 and Acrobot-v1. The objective in CartPole-v0 is to balance a pole on a cart. Each step generates a positive reward signal while receiving negative rewards if the pole falls below an angle threshold. The agent can control the direction of the cart at every time step. The Acrobot-v1 has a similar aim to control the arm to hit the ceiling in a minimal number of

**Fig. 2.** The experiment test-bed contains the following environments; CartPole-v0, Acrobot-v1, CrossingS9N3-v0, and FourRooms-v0

time steps. The agent receives negative rewards until it reaches the ceiling. The CartPole-v0 and Acrobot-v1 environments origins from [4] while CrossingS9N3-v0 and the FourRooms-v0 origins from [6].[2]

### 5.2   Hyperparameters

During the experiments, we found DVQN to be challenging to tune. Initially, the algorithm used ReLU as activation but was discarded due to vanishing gradients resulting in divergence for both policy and reconstruction objectives. By using ELU, we found the algorithm to be significantly more stable during the experiments, and it additionally did not diverge if training continued after convergence. We will explore the underlying cause of our future work. Table 1 shows the hyperparameters used in our experiments where most of the parameters are adopted from prior work. Recognize that the DVQN algorithm does not use $\epsilon$-greedy methods for exploration. The reason for this is that random sampling is done during training in the variational autoencoder part of the architecture. In general, the algorithm tuning works well across all of the tested domains, and better results can likely be achieved with extended hyperparameter searches. For DVQN, we consider such tuning for our continued work on DVQN using options, see Section 6.

---

[2] A community-based scoreboard can be found at `https://github.com/openai/gym/wiki/Leaderboard`.

| Algorithm | DQN | DDQN | Rainbow | DVQN (ours) |
|---|---|---|---|---|
| Optimiser | | Adam | | RMSProp |
| Learning Rate | | 0.003 | | 0.000025 |
| Activation | | ReLU | | ELU |
| Batch Size | | 32 | | 128 |
| Replay Memory | | 1m | | |
| Epsilon Start | | 1.0 | | N/A |
| Epsilon End | | 0.01 | | N/A |
| Epsilon Decay | | 0.001 (Linear) | | N/A |
| Gamma | | 0.95 | | |
| Q-Loss | | Huber | | MSE |

**Table 1.** Algorithm and hyperparameters used in the experiments. For the Rainbow algorithm, we used the same hyperparameters described in [10]. The DDQN had target weight updates every 32k frames.
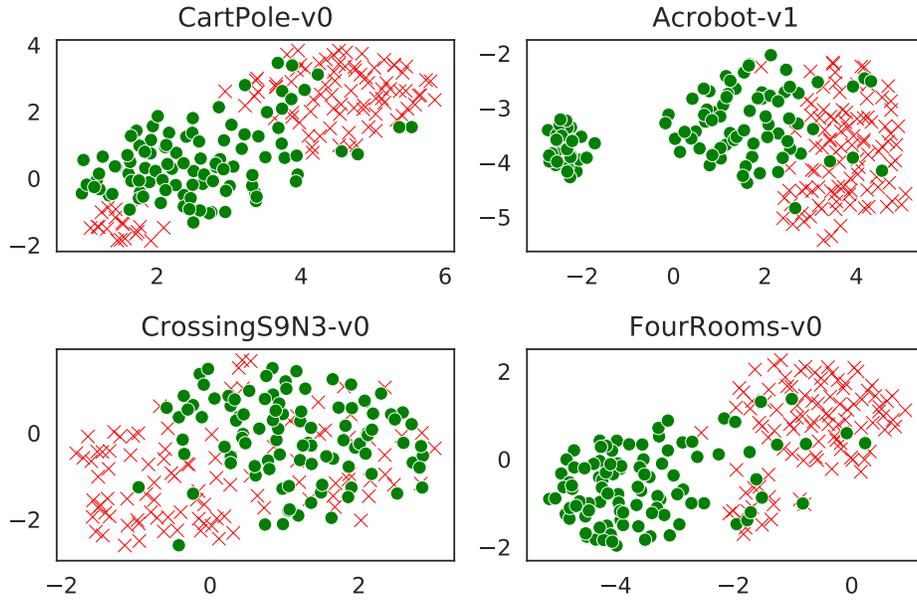
### 5.3   Latent-Space evaluation

An attractive property of our model is that the latent-space is a Gaussian distribution. As seen in Figure 3, the DVQN algorithm can produce clustered latent-spaces for all tested environments. For example, in the CartPole-v0 environment, there are three clusters where two of them represent possible terminal states and one that represents states that give a reward. To fully utilize the capabilities of DVQN, the latent-space can be used to generate options for each cluster to promote different behavior for every region of the state-space.

Figure 4 illustrates the visualization of the latent-space representation in CartPole-v0. We find that each cluster represents a specific position and angle of the pole. The latent-space interpolates between these state variations, which explains its shape. Although the clusters are not perfect, it is trivial to construct separate classification for each cluster with high precision, and this way automatically construct initiation and termination signals for options.
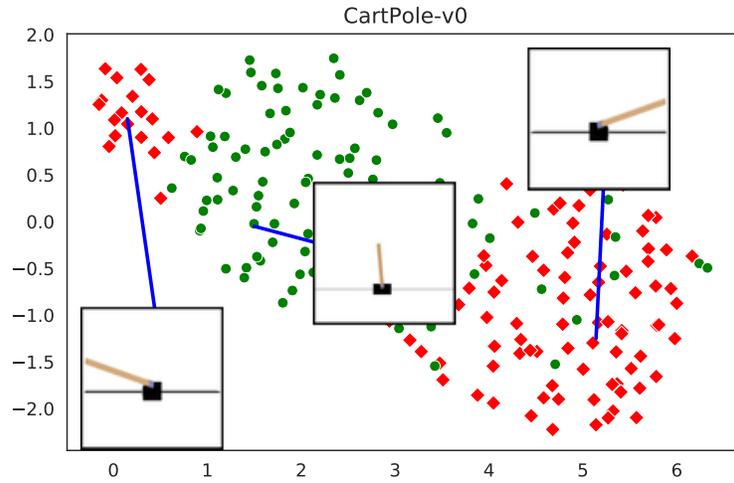
### 5.4   Performance evaluation

Figure 5 illustrates a comparison of performance between state-of-the-art Q-Learning algorithms and the proposed DVQN algorithm. The performance measurement is the mean of 100 trials over 1500 episodes for CartPole-v0, CrossingS9N3-v0, FourRooms-v0, and 3000 episodes for Acrobot-v1. The performance is measured in accumulated rewards and is therefore negative for environments where each time step yields a negative reward.
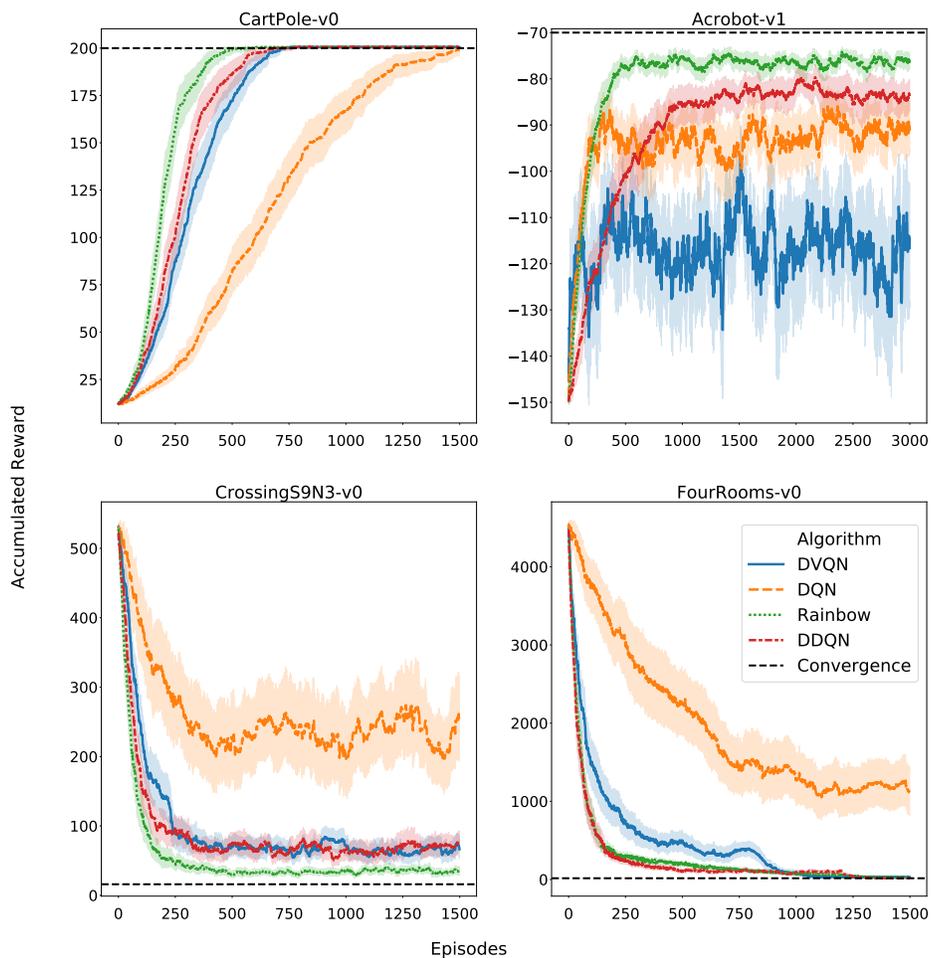
The DVQN algorithm performs better than DQN and shows comparable performance to DDQN and Rainbow. DVQN is not able to find a good policy in the Acrobot-v1 environment but successfully learns a good visual representation of the latent space. In general, the DVQN algorithm is significantly harder to train because it requires the algorithm to find a good policy within a Gaussian distribution. We found this to work well in most cases, but it required fine-tuning

**Fig. 3.** The learned latent space for all of the tested environments. DVQN successfully trivialise the selection of options as seen in the well-separated state clusters. The circular points illustrate states with positive reward while cross illustrates negative rewards.



**Fig. 4.** The relationship between states and the latent-space for the CartPole-v0 environment. DVQN can separate each angle, left, middle and right into separable clusters, which are especially useful in option-based reinforcement learning. Additionally, the visualization of the latent space that the Q-head uses to sample actions is trivial to interpret.

**Fig. 5.** The accumulative sum of rewards of the DVQN compared to other Q-Learning based methods in the experimental environments. Our algorithm performs better than DQN from [14], and shows comparable results to DDQN from [9] and Rainbow from [10]. We define an episode threshold for each of the environments (x-axis) and accumulate the agent rewards as the performance metric for CartPole-v0 and Acrobot-v1. The scoring metric in CrossingS9N3-v0 and the FourRooms-v0 is based on how many steps the agent used to reach the goal state.

of hyperparameters. The algorithm is also slower to converge, but we were able to improve training stability by increasing the batch-size and decreasing the learning-rate.

## 6    Conclusion and Future Work

This paper introduces the deep variational Q-network (DVQN), a novel algorithm for learning policies from a generative latent-space distribution. The learned latent-space is particularly useful for clustering states that are close to each other for **discovering options automatically**. In the tested environments, the DVQN algorithm can achieve **comparable performance** to traditional deep Q-networks. DVQN does not provide the same training stability and is significantly harder to fine-tune than traditional deep Q-learning algorithms. For instance, network capacity is increased. As a result of this, the algorithm takes longer to train, and during the experiments, only the RMSprop optimizer [15] with a small step size was able to provide convergence. Additionally, the exponential linear units from [7] had a positive effect on stability. On the positive side, the DVQN contributes a **novel approach for options discovery** in hierarchical reinforcement learning algorithms.

The combination of VAE and reinforcement learning algorithms has interesting properties. Under **optimal conditions**, the latent-space should, in most cases follow a true Gaussian distribution where policy evaluations always provide optimal state-action values, since this is the built-in properties of the latent space in any VAE. The difference between traditional deep Q-networks and DVQN primarily lies in the elimination of a sparse and unstructured latent-space. In deep Q-Networks, optimization does not provide a latent-space structure that reflects a short distance between states but rather a distance between Q-values [14]. By using KL-regularization from VAE, low state-to-state is encouraged. Another benefit of VAE is that we sample from a Gaussian distribution to learn $\mu$ and $\sigma$, which is especially satisfying for algorithms with off-policy sampling and therefore eliminates the need for ($\epsilon$-greedy) random sampling.

In the continued work, we wish to do a thorough analysis of the algorithm to justify its behavior and properties better. A better understanding of the Gaussian distributed latent-space is particularly appealing because it would enable better labeling schemes for clustering, or perhaps fully automated labeling. Finally, we plan to extend the algorithm from model-free behavior to hierarchical RL with options. The work of this contribution shows that it is feasible to produce organized latent-spaces that could provide meaningful options, and the hope is that this will result in state-of-the-art performance in a variety of tasks in RL.

## References

1. Annasamy, R.M., Sycara, K.: Towards Better Interpretability in Deep Q-Networks. Proceedings, The Thirty-Third AAAI Conference on Artificial Intelligence (sep 2018), http://arxiv.org/abs/1809.05630

2. Arulkumaran, K., Deisenroth, M.P., Brundage, M., Bharath, A.A.: Deep reinforcement learning: A brief survey. IEEE Signal Processing Magazine **34**(6), 26–38 (2017). https://doi.org/10.1109/MSP.2017.2743240

3. Barto, A., Mahadevan, S., Lazaric, A.: Recent Advances in Hierarchical Reinforcement Learning. Tech. rep., PIGML Seminar-AirLab (2003)

4. Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W.: OpenAI Gym. arxiv preprint arXiv:1606.01540 (jun 2016), `http://arxiv.org/abs/1606.01540`

5. Chen, W., Zhang, M., Zhang, Y., Duan, X.: Exploiting meta features for dependency parsing and part-of-speech tagging. Artificial Intelligence **230**, 173–191 (sep 2016). https://doi.org/10.1016/j.artint.2015.09.002, `http://arxiv.org/abs/1509.06461`

6. Chevalier-Boisvert, M., Willems, L., Pal, S.: Minimalistic Gridworld Environment for OpenAI Gym. \url{https://github.com/maximecb/gym-minigrid} (2018)

7. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs). The International Conference on Learning Representations 16 (nov 2015), `http://arxiv.org/abs/1511.07289`

8. Ha, D., Schmidhuber, J.: Recurrent World Models Facilitate Policy Evolution. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) Advances in Neural Information Processing Systems 31, pp. 2450–2462. Curran Associates, Inc., Montréal, CA (sep 2018), `http://papers.nips.cc/paper/7512-recurrent-world-models-facilitate-policy-evolution.pdf`

9. van Hasselt, H., Guez, A., Silver, D.: Deep Reinforcement Learning with Double Q-learning. Proceedings, The Thirtieth AAAI Conference on Artificial Intelligence p. 13 (sep 2015), `http://arxiv.org/abs/1509.06461`

10. Hessel, M., Modayil, J., van Hasselt, H., Schaul, T., Ostrovski, G., Dabney, W., Horgan, D., Piot, B., Azar, M., Silver, D.: Rainbow: Combining Improvements in Deep Reinforcement Learning. In: Proc. 32nd Conference on Artificial Intelligence, AAAI'18. pp. 3215–3222. AAAI Press, New Orleans, Louisiana USA (oct 2018), `https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/download/17204/16680`

11. Huang, S., Su, H., Zhu, J., Chen, T.: SVQN: Sequential Variational Soft Q-Learning Networks. In: International Conference on Learning Representations (2020), `https://openreview.net/forum?id=r1xPh2VtPB`

12. Kingma, D.P., Welling, M.: Auto-Encoding Variational Bayes. Proceedings of the 2nd International Conference on Learning Representations (dec 2013). https://doi.org/10.1051/0004-6361/201527329, `http://arxiv.org/abs/1312.6114`

13. Levine, S., Finn, C., Darrell, T., Abbeel, P.: End-to-end training of deep visuomotor policies. Journal of Machine Learning Research **17**(1), 1334–1373 (2016), `http://www.jmlr.org/papers/volume17/15-522/15-522.pdf`

14. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning (dec 2015). https://doi.org/10.1038/nature14236, `http://arxiv.org/abs/1312.5602`

15. Nair, V., Hinton, G.E.: Rectified Linear Units Improve Restricted Boltzmann Machines. Proceedings of the 27 th International Conference on Machine Learning p. 8 (2010)

16. Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., Hubert, T., Baker, L., Lai, M., Bolton, A., Chen, Y., Lillicrap, T., Hui, F., Sifre, L., Van Den Driessche, G., Graepel, T., Hassabis, D.: Mastering the game of Go without human knowledge. Nature (2017). https://doi.org/10.1038/nature24270
17. Stolle, M.: Automated Discovery of Options in Reinforcement Learning. Ph.D. thesis, McGill University (2004)
18. Sutton, R.S., Precup, D., Singh, S.: Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning. Artificial Intelligence **112**(1-2), 181–211 (1999)
19. Tang, Y., Kucukelbir, A.: Variational Deep Q Network. In: Advances in Neural Information Processing Systems 30. Long Beach, CA, USA (nov 2017), `http://arxiv.org/abs/1711.11225`
20. Wang, J., Gou, L., Shen, H.W., Yang, H.: DQNViz: A Visual Analytics Approach to Understand Deep Q-Networks. IEEE Transactions on Visualization and Computer Graphics **25**(1), 288–298 (jan 2019). https://doi.org/10.1109/TVCG.2018.2864504
21. Younes, H.L.S., Simmons, R.G.: Solving Generalized Semi-Markov Decision Processes using Continuous Phase-Type Distributions. Proceedings, The Ninth AAAI Conference on Artificial Intelligence (2004), `www.aaai.org`
22. Zhang, C., Patras, P., Haddadi, H.: Deep Learning in Mobile and Wireless Networking: A Survey. IEEE Communications Surveys & Tutorials (mar 2018), `http://arxiv.org/abs/1803.04311`
23. Zheng, Y., Tan, H., Tang, B., Zhou, H., Jiang, Z., Zheng, Y., Tan, H., Tang, B., Zhou, H.: Variational Deep Embedding: A Generative Approach to Clustering. International Joint Conference on Artificial Intelligence 17 p. 8 (2017), `http://arxiv.org/abs/1611.05148`