

Ising Networks for Deep Hierarchical Reinforcement Learning

1 Introduction

2 Related Work

2.1 Hierarchical Reinforcement Learning

2.1.1 Temporal Abstraction

Various methods have devised hierarchical techniques to abstract the temporal correlation in MDPs [16, 17, 30, 35]. The Hierarchical Deep Q-Network (hDQN) [16] proposes an integration of temporal abstraction and intrinsic motivation in order to solve long-horizon problems. hDQN demonstrates improved long-horizon optimality at the cost of manual feature engineering of the states and reward function. This is addressed by making use of deep successor networks [17] which extract relevant goal embeddings for the agent and direct its behavior towards the goal. A more sophisticated temporal scheme consists of learning abstraction of hierarchies based on expert demonstrations and regularized latent space constraints [38]. However, sub-policies rely on memory-based controllers for learning fine-grained abstraction and often cripple the lower levels of the hierarchy once the task has been accomplished. This leads to limited transfer of skills between policies and hence, results in sub-optimal convergence [7]. Transfer of skills between different policies operating in sub-MDPs is extensively studied using exploration [15]. [25] presents MAVEN, a hierarchical architecture for variational exploration in the case of multi-agent settings. Policies carry out exploration in the latent space using variational inference. MAVEN demonstrates improved coordination between agents in the long-horizon. However, policy behavior distributed across agents which forces sub-policies to rely on other agent policies for optimal convergence. [4] proposes the trajectory autoencoder which enables the agent to exhibit self-consistent behavior and tackle sparse rewards in the long-horizon. Long-horizon suitability can further be improved by jointly training all policies and utilize the skills learned individually [20]. Such kind of sub-policy adaptation stabilizes the behavior of agent and provides efficient training of hierarchies at all levels.

Abstraction may be carried out among states [2] or using a different medium such as language [14] which allows the agent to solve temporally-extended tasks in a diverse manner. [5] presents the MAXQ algorithm which proposes an efficient method for training state-based abstractions in the form of hierarchies. Q-values corresponding to each lower-level task are aggregated at the higher-level policy nodes to yield a modular framework for abstraction. While the MAXQ demonstrates suitability for low-dimensional tasks, its performance cannot be assessed in high-dimensional state and action spaces due to the larger number of abstractions produced by the framework. [1] tackles the problem of high-dimensional and long-horizon learning by presenting a theoretical perspective of abstracting states using transitive behavior which helps in optimal computations. Such techniques have motivated active research in the form of modified abstraction algorithms for performance improvement and schemes leading to richer abstract representations [22].

2.1.2 Option-based Hierarchies

Hierarchies can be constructed on the basis of skills or possible options of policies which the agent may choose from during its interactions in the environment [35, 3]. [3] presents the Options-Critic Architecture which is a collection of various policy options. Agent can select from the collection of sub-policies at various timesteps to yield a policy which is robust and optimal in the long-horizon. The provision of options enables the agent to acquire a wide variety of skills during exploration and obtain a policy which makes best use of these skills. Furthermore, the architecture incorporates learning of termination conditions of a particular option which is essential for switching option policies [11]. The architecture is extended to the Double Actor-Critic architecture which comprises of two parallel MDPs [40]. A dual architecture enables the learning of intra-option policies with termination conditions of each option. While, option-based learning provides diversity in skills and control of temporal extension to the agent, frequent fluctuations between options results in instability in updates and hence, sample-inefficient learning [21]. [21] addresses the instability in training by combining the options framework with off-policy maximum entropy reinforcement learning [10, 9]. Maximization of entropy allows the agent to effectively explore all options and consistently select the optimal sub-policies based on selected options. However, usage of different skills and compactness remains an open problem in the case of real-world tasks consisting of large action spaces. [28] addresses the problem of compact options by improvising communication between policies using binary vectors. Modular policies using binary vectors enable richer communication and mixing of skills. Techniques related to the relative optimization of policy can be extended using advantage-weighted importance sampling [27] in options based on option-value functions.

2.1.3 Hierarchical Control

Hierarchical Reinforcement Learning borrows from multiple approaches in control [29] and adaptive policy optimization [20]. Multiplicative Compositional Policies (MCPs) [29] present the composable nature of hierarchies which arises from coordination of multiple control-based skills. MCPs improve the hierarchical coordination between sub-policies which results in adaptive motor skills in the case of robotic control such as locomotion and manipulation. While MPC demonstrates the practical usage of hierarchical reinforcement learning motivated by adaptive control, [13] extends the hierarchical framework to Quadruped Locomotion. The lower-level policy uses latent commands for robot actuator control and the higher-level policy operates at different time intervals. Abstraction of control between lower and higher-level policies allows the agent to acquire a diverse set of task-related skills. Moreover, policies can be efficiently adapted on a different task from the same domain. Additionally, diversity can be achieved in a scalable manner by training policies levelwise [32]. However, levelwise training of sub-policies becomes intractable in cases of a large number of hierarchies [8]. [8] tackles the problem of a large number of hierarchies and levelwise training by making use of latent space policies. Each latent space policy is trained using the maximum entropy learning framework in order to directly solve the task. Hierarchies trained in latent space demonstrate suitable initializations on complex control tasks which enable the agent to maintain consistent behavior. Consistency can also be achieved using goal embeddings which motivate goal-directed behavior in the case of navigation tasks [34]. [37] additionally improves consistency in hierarchical agents by making use of information theoretic regularizations which pose a reward penalty on the policy distribution. Hierarchical control in reinforcement learning has been successful in tackling knowledge transfer as well as data-efficiency [39]. Various levels of the hierarchy make use of inductive-biases for sharing off-policy data. Sharing knowledge across policies in a compositional manner results in positive transfer of policies. Such techniques are additionally used to solve numerically challenging control problems consisting of unstructured environments [23].

2.2 Ising Models

2.2.1 Energy-based Reinforcement Learning

Recent work on energy-based learning [19] has demonstrated significant success in the arena of reinforcement learning [33, 18]. [31] highlights the use of energy-based policies capable of learning large action spaces by making use of factored approximations. Such approximations can be thought of as inference mechanisms [24] using energy-based distributions such as the boltzmann distribution [36]. These have given rise to energy-based extensions of models with different architectures such as actor-critic frameworks [12]. Another suitable method for training energy-based policies in reinforcement learning is by making use of the Helmholtz free energy [9]. [?] presents the suitability of improved helmholtz energy operator under various reinforcement learning scenarios with theoretical

guarantees on robustness and improved inference in comparison to boltzmann energy. Applications of the energy operator include surprise minimization in the case of multi-agent learning [?]. [?] highlights the practical use of energy in multi-agent learning across agents which improves the joint policy to tackle environments consisting of high variance. The use of energy-based policies can further be generalized to inverse reinforcement learning [6] wherein the maximum entropy framework can be realized as a special case of energy-based models in the generative setting. Although energy plays a key role the training dynamics of reinforcement learning and optimization of the agent to an optimal policy, it does little to improve the computational framework of learning hierarchical relationships in agent’s policy [26]. This indicates the requirement for a hierarchical framework capable of mapping relationships between sub-policies in the hierarchy.

2.2.2 Ising Model Learning

3 The Ising Model

4 Intuition for Spin Values

5 Implementation Details

References

- [1] D. Abel, D. Arumugam, L. Lehnert, and M. Littman. State abstractions for lifelong reinforcement learning. In *International Conference on Machine Learning*, 2018.
- [2] D. Andre and S. J. Russell. State abstraction for programmable reinforcement learning agents. In *AAAI/IAAI*, 2002.
- [3] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture, 2016.
- [4] J. D. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine. Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings. *arXiv preprint arXiv:1806.02813*, 2018.
- [5] T. G. Dietterich. State abstraction in maxq hierarchical reinforcement learning. In *Advances in Neural Information Processing Systems*, 2000.
- [6] C. Finn, P. Christiano, P. Abbeel, and S. Levine. A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models, 2016.
- [7] C. Florensa, Y. Duan, and P. Abbeel. Stochastic neural networks for hierarchical reinforcement learning, 2017.

- [8] T. Haarnoja, K. Hartikainen, P. Abbeel, and S. Levine. Latent space policies for hierarchical reinforcement learning, 2018.
- [9] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies, 2017.
- [10] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor, 2018.
- [11] A. Harutyunyan, W. Dabney, D. Borsa, N. Heess, R. Munos, and D. Precup. The termination critic, 2019.
- [12] N. Heess, D. Silver, and Y. W. Teh. Actor-critic reinforcement learning with energy-based policies. *Proceedings of Machine Learning Research*, 2013.
- [13] D. Jain, A. Iscen, and K. Caluwaerts. Hierarchical reinforcement learning for quadruped locomotion. *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.
- [14] Y. Jiang, S. S. Gu, K. P. Murphy, and C. Finn. Language as an abstraction for hierarchical deep reinforcement learning. In *Advances in Neural Information Processing Systems 32*. 2019.
- [15] N. K. Jong, T. Hester, and P. Stone. The utility of temporal abstraction in reinforcement learning. In *AAMAS (1)*, 2008.
- [16] T. D. Kulkarni, K. R. Narasimhan, A. Saeedi, and J. B. Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation, 2016.
- [17] T. D. Kulkarni, A. Saeedi, S. Gautam, and S. J. Gershman. Deep successor reinforcement learning, 2016.
- [18] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas. Reinforcement learning with augmented data, 2020.
- [19] Y. LeCun, S. Chopra, R. Hadsell, F. J. Huang, and et al. A tutorial on energy-based learning. In *PREDICTING STRUCTURED DATA*, 2006.
- [20] A. Li, C. Florensa, I. Clavera, and P. Abbeel. Sub-policy adaptation for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2020.
- [21] C. Li, X. Ma, C. Zhang, J. Yang, L. Xia, and Q. Zhao. Soac: The soft option actor-critic architecture, 2020.
- [22] S. Lin and R. Wright. Evolutionary tile coding: An automated state abstraction algorithm for reinforcement learning. In *Proceedings of the 8th AAAI Conference on Abstraction, Reformulation, and Approximation*, 2010.

- [23] Q. Ma, S. Ge, D. He, D. Thaker, and I. Drori. Combinatorial optimization by graph pointer networks and hierarchical reinforcement learning, 2019.
- [24] D. J. C. MacKay. *Information Theory, Inference & Learning Algorithms*. Cambridge University Press, 2002.
- [25] A. Mahajan, T. Rashid, M. Samvelyan, and S. Whiteson. Maven: Multi-agent variational exploration. In *Advances in Neural Information Processing Systems 32*. 2019.
- [26] O. Nachum, S. Gu, H. Lee, and S. Levine. Near-optimal representation learning for hierarchical reinforcement learning, 2018.
- [27] T. Osa, V. Tangkaratt, and M. Sugiyama. Hierarchical reinforcement learning via advantage-weighted information maximization, 2019.
- [28] A. Pashevich, D. Hafner, J. Davidson, R. Sukthankar, and C. Schmid. Modulated policy hierarchies, 2018.
- [29] X. B. Peng, M. Chang, G. Zhang, P. Abbeel, and S. Levine. Mcp: Learning composable hierarchical control with multiplicative compositional policies, 2019.
- [30] D. Precup. Temporal abstraction in reinforcement learning, 2000.
- [31] B. Sallans and G. E. Hinton. Reinforcement learning with factored states and actions. *The Journal of Machine Learning Research*, 2004.
- [32] Y. Song, J. Wang, T. Lukasiewicz, Z. Xu, and M. Xu. Diversity-driven extensible hierarchical reinforcement learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33, 2019.
- [33] A. Srinivas, M. Laskin, and P. Abbeel. Curl: Contrastive unsupervised representations for reinforcement learning, 2020.
- [34] S. Sukhbaatar, E. Denton, A. Szlam, and R. Fergus. Learning goal embeddings via self-play for hierarchical reinforcement learning, 2018.
- [35] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. 2018.
- [36] Y. W. Teh, M. Welling, S. Osindero, and G. E. Hinton. Energy-based models for sparse overcomplete representations. *The Journal of Machine Learning Research*, 2003.
- [37] D. Tirumala, H. Noh, A. Galashov, L. Hasenclever, A. Ahuja, G. Wayne, R. Pascanu, Y. W. Teh, and N. Heess. Exploiting hierarchy for learning and transfer in kl-regularized rl, 2019.
- [38] W. Wang, Y. Hu, and S. Scherer. Learning temporal abstraction with information-theoretic constraints for hierarchical reinforcement learning, 2020.

- [39] M. Wulfmeier, A. Abdolmaleki, R. Hafner, J. Tobias Springenberg, M. Neunert, N. Siegel, T. Hertweck, T. Lampe, N. Heess, and M. Riedmiller. Compositional transfer in hierarchical reinforcement learning. *Robotics: Science and Systems XVI*, 2020.
- [40] S. Zhang and S. Whiteson. Dac: The double actor-critic architecture for learning options, 2019.