

# IP Summary - Reinforcement Learning

Aman Mehra

December 23, 2020

This project began with the aim of surveying and understanding the existing literature in Reinforcement Learning. In the first few weeks, I performed a literature survey of existing DQN architectures that are employed in the field. The four major architectures that I surveyed were - DQN [5], Double DQN [3], Dueling Networks [10] and Rainbow [4]. Minh *et al* [5] enabled deep neural networks to be employed as functional approximators for Q learning by handling delayed, sparse and noisy rewards and breaking data correlation through replay memories. The subsequent papers incrementally improved this architecture. Double DQN handled the over estimation bias in Q-learning due the  $\max()$  function by using a second network to obtain bootstrap values for the Q value target. Dueling networks uses two streams to disentangle the value function and action advantage function with the aim of learning good state estimates without needing to observe all possible actions from this state. Rainbow combines the recent advances in RL architectures, experience replay and rewards to achieve state of the art performance on the Arcade Learning benchmark. Alongside, I also explored the area of intrinsic motivation through the paper on curiosity driven exploration [7].

Subsequently, the focus shifted towards understanding replay memories. Initial works in this domain such as PER [8], CER [12] and PSER [1] were first studied. These focused on TD error based prioritization of transitions in the replay memories. Newer methods follow a different approach by introducing the notion of on-policy-ness. These works try to maintain more on policy transitions with respect to the current agent policy in the replay memory. The two major works that were analyzed were ERO [11] and ReF-ER [6]. ERO learns a replay policy that chooses samples that maximize the improvement in the cumulative reward of the agent. This approach surpasses TD error based prioritization methods while also contradicting the older paradigm by picking low TD-error transitions instead of high TD-error transitions. This finding puts the notion of on-policy transition selection to the fore. ReF-ER presents an alternate to pick on-policy transitions. It solves the problem by using the importance sampling ratio to classify transitions as near[on] or far[off] policy transitions and back-propagating TD error of only near policy transitions. It also uses KL divergence loss to smoothen the change of the current target policy with respect to an older target policy. Effectively, this paper presents a scheme to smoothen the change in policy of an agent in the later stages of learning. Additionally, an orthogonal

approach of Importance Resampling [9] was analyzed wherein transitions are sampled with a probability proportional to their IS ratio to reduce the overall variance of the update.

In the remaining period we analyzed the findings of a recent paper[2] that studied the interplay between replay memory parameters and different aspects of the RL pipeline to determine what parts enabled improved returns with larger replay memories. Their findings were that n-step returns were the key to seeing improved returns with larger replay memories as the larger memory compensated for the added variance of the update target. Since then we have been trying to reproduce and verify the results of this paper’s experiments. The subsequent line of work is to investigate whether this hypothesis holds for low variance multi step returns like tree back up and n-step average DQN bootstraps. Concretely, we want to investigate the hypothesis that - larger replay memories reduce the variance of the target  $E[G_t]$  and thus help mitigate the added variance of the n-step return. We aim to study the trends of this variance for different returns to gain a deeper understanding of the phenomenon.

## References

- [1] BRITTAİN, M., BERTRAM, J. R., YANG, X., AND WEI, P. Prioritized sequence experience replay. *CoRR* (2019).
- [2] FEDUS, W., RAMACHANDRAN, P., AGARWAL, R., BENGIO, Y., LAROCHELLE, H., ROWLAND, M., AND DABNEY, W. Revisiting fundamentals of experience replay. In *International Conference on Machine Learning* (2020).
- [3] HASSELT, H. V., GUEZ, A., AND SILVER, D. Deep reinforcement learning with double q-learning. *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*.
- [4] HESSEL, M., MODAYIL, J., VAN HASSELT, H., SCHAUL, T., OSTROVSKI, G., DABNEY, W., HORGAN, D., PIOT, B., AZAR, M. G., AND SILVER, D. Rainbow: Combining improvements in deep reinforcement learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence* (2018).
- [5] 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., ANTONOGLU, I., KING, H., KUMARAN, D., WIERSTRA, D., LEGG, S., AND HASSABIS, D. Human-level control through deep reinforcement learning. *Nature* (2015).
- [6] NOVATI, G., AND KOUMOUTSAKOS, P. Remember and forget for experience replay. In *Proceedings of the 36th International Conference on Machine Learning* (2019).

- [7] PATHAK, D., AGRAWAL, P., EFROS, A. A., AND DARRELL, T. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the 34th International Conference on Machine Learning* (2017).
- [8] SCHAUL, T., QUAN, J., ANTONOGLOU, I., AND SILVER, D. Prioritized experience replay. In *International Conference on Learning Representations* (2016).
- [9] SCHLEGEL, M., CHUNG, W., GRAVES, D., QIAN, J., AND WHITE, M. Importance resampling for off-policy prediction. In *Advances in Neural Information Processing Systems* (2019).
- [10] WANG, Z., SCHAUL, T., HESSEL, M., VAN HASSELT, H., LANCTOT, M., AND DE FREITAS, N. Dueling network architectures for deep reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning* (2016).
- [11] ZHA, D., LAI, K.-H., ZHOU, K., AND HU, X. Experience replay optimization. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19* (2019).
- [12] ZHANG, S., AND SUTTON, R. S. A deeper look at experience replay. In *Symposium on Deep Reinforcement Learning at the 31st Conference on Neural Information Processing Systems* (2017).