

Tactical Optimism and Pessimism for Deep Reinforcement Learning

Ted Moskovitz, Jack Parker-Holder, Aldo Pacchiano,
Michael Arbel, Michael I. Jordan

NeurIPS 2021



**Gatsby Computational
Neuroscience Unit**

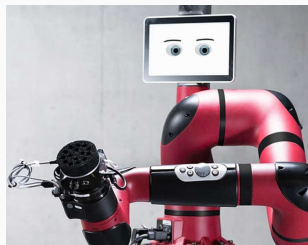
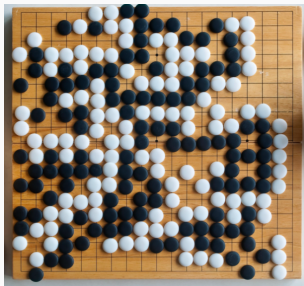
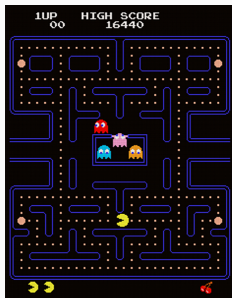


UNIVERSITY OF
OXFORD



Berkeley
UNIVERSITY OF CALIFORNIA

Motivation

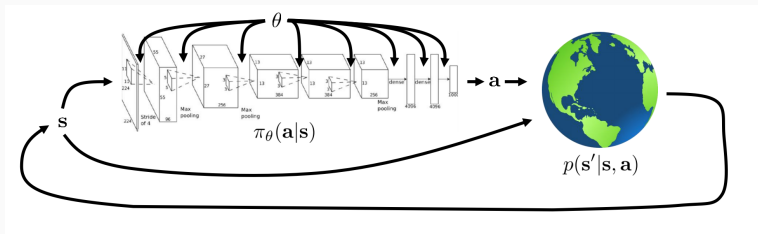


- Approximate value functions are key to the success of deep RL
- Optimism vs. pessimism: both have drawbacks + advantages
- How can we figure out which is best?

Summary + Contributions

- Demonstrate that the efficacy of optimism varies both across environments and over the course of training
- Introduce a novel framework for value estimation, *Tactical Optimism and Pessimism* (TOP)
- Adaptively updates its degree of optimism by modeling the choice as a multi-armed bandit problem
- Augmenting popular algorithms with TOP leads to state-of-the-art results

RL Background



- Assume an agent is acting in an MDP $(S; A; r; p; \gamma)$
- Running the policy in an episodic/finite horizon task of length T produces trajectories $\tau = (s_1; a_1; r_1; \dots; s_T; a_T; r_T)$
- Maximize $J(\pi) = E_{\pi} [Z(\pi)] = E_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right]$

The Actor-Critic Framework

- Actor, π : a deterministic policy
- Critic, Q : evaluate actor,

$$Q(s; a) := E [Z_t | s_t = s; a_t = a] \quad (0.1)$$

- How to train?
- Actor: $J(\pi) = E \sum_{t=0}^{\infty} \gamma^t r_{a_t} Q(s_t; a_t) |_{a_t = \pi(s_t)}$
- Critic: given a transition $(s_t; a_t; r_t; s_{t+1})$,

$$\nabla_{\theta} \frac{1}{2} \| y_t - Q(s_t; a_t) \|^2 \quad (0.2)$$

$y_t = r_t + \gamma Q(s_{t+1}; a_{t+1})$
 $Z_t = \sum_{i=t}^{\infty} \gamma^i (r_{s_i, a_i} + \gamma Q(s_{i+1}; a_{i+1}) - Q(s_i; a_i))$

where $y_t = r_t + \gamma Q(s_{t+1}; a_{t+1})$

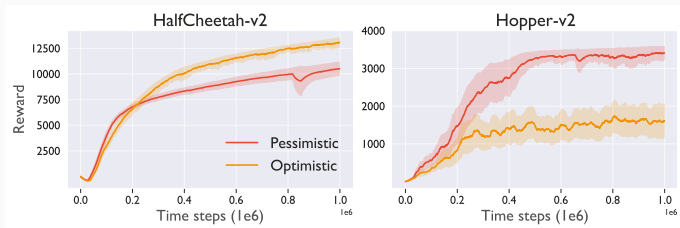
- Distributional RL: represent *distribution* of Z , not just mean [1]

Optimism vs. Pessimism

- **Problem:** Function approx. biases critic towards **overestimation** [5]
- **Solution:** build a **pessimistic** target using 2 critics [3]:

$$y_t = r_t + \min_{i \in \{1, 2\}} Q_i(s_{t+1}; (s) +); \quad \text{clip}(N(0; s^2); -c; c) \quad (0.3)$$

- **Problem:** pessimistic critics can result in **underexploration** [2]
- **Solution:** use an **optimistic** upper-bound on the value [2]
- Confusing...



TOP: Tactical Optimism and Pessimism

- Distinguish between and represent two types of uncertainty:
- **Aleatoric uncertainty:** noise inherent to the world/task
- $q_{Z^*}^{(k)}$ represent using *distributional value estimation* $Z^* \sim Z(s; a)$
- **Epistemic uncertainty:** noise due to lack of knowledge about the world
- $q_{Z^{\wedge}}^{(k)}$ represent using an *ensemble* of $k = 1; \dots; K$ critics
- We use these estimates to construct a *belief distribution* $Z^{\wedge}(s; a)$:

$$q_{Z^{\wedge}}^{(k)}(s; a) = q_{Z^*}^{(k)}(s; a) + q_{Z^{\wedge}}^{(k)}(s; a) \quad (0.4)$$

- β then determines the degree of optimism

TOP: Tactical Optimism and Pessimism

- Q: How to choose ?
- A: Evaluate by its **effect on performance!**
- Model choice as a **bandit problem**:
 - Choose from $\{\beta_d\}_{d=1}^D$ by sampling a decision $d_m \in \{1, \dots, D\}$ for episode m
 - $d_m \sim p_m(\cdot)$, where $p_m(d) \propto \exp(w_m(d))$
- Update arm weightings:

$$w_{m+1}(d) = \begin{cases} w_m(d) + \frac{f_m}{p_m(d)} & \text{if } d = d_m \\ w_m(d) & \text{otherwise;} \end{cases} \quad (0.5)$$

- Feedback f_m is **change in performance**:

$$f_m = R_m - R_{m-1} \quad (0.6)$$

Experiments: State-based Control

TOP + TD3 [3]:



Experiments: State-based Control

TOP + TD3 [3]:



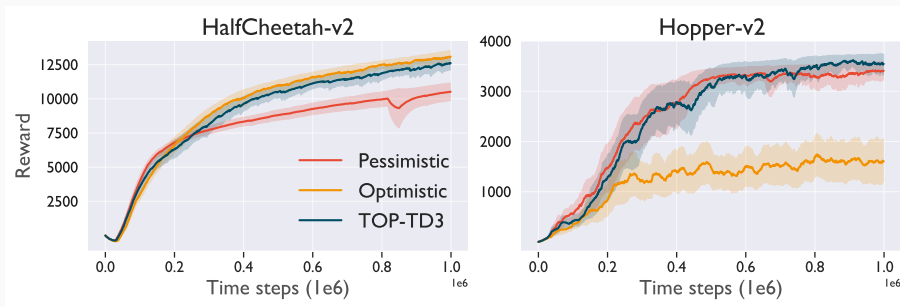
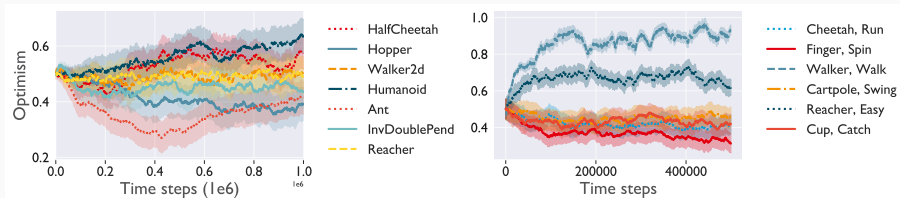
Experiments: Pixel-based Control

TOP + RAD [4]:

Experiments: Pixel-based Control

TOP + RAD [4]:

Experiments: The Impact of Adaptive Optimism



Summary

- It's difficult to know set the correct degree of optimism
- TOP is an adaptive, uncertainty-based method which does it for you
- TOP boosts SOTA performance on state- and pixel-based control
- Adding TOP requires only 10 lines of Python code

Thank you!

- [1] M. G. Bellemare, W. Dabney, and R. Munos. A distributional perspective on reinforcement learning. *arXiv preprint arXiv:1704.06768*, 2017.
- [2] K. Ciosek, Q. Vuong, R. Loftin, and K. Hofmann. Better exploration with optimistic actor-critic, 2019.
- [3] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. *CoRR*, abs/1802.09477, 2018.
- [4] M. Laskin, K. Lee, A. Stooke, L. Pinto, P. Abbeel, and A. Srinivas. Reinforcement learning with augmented data. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19884–19895. Curran Associates, Inc., 2020.
- [5] S. Thrun and A. Schwartz. Issues in using function approximation for reinforcement learning. In *In Proceedings of the Fourth Connectionist Models Summer School*. Erlbaum, 1993.