

Deep Policy Gradient Methods Without Batch Updates, Target Networks, or Replay Buffers

Gautham Vasan, Mohamed Elsayed, Alireza Azimi, Jiamin He, Fahim Shahriar, Colin Bellinger, Martha White & A. Rupam Mahmood

Deep Learning: Classics and Trends, ML Collective Feb 7th 2025











I'm interested in agents that can learn on the fly, in real-time, by interacting with their environment

Fun with Squirrels

- Real-time learning
- Learn through trial-and-error interactions with the environment
- React in milliseconds to achieve their goal





Real-Time Adaptation



Application of an AR Headset



Mars Perseverance Rover

- system frozen during deployment
- changes in dynamic environments
 - Learning must occur on-device!







Apple Watch

iRobot Roomba

In current AI systems, training and deployment are distinct phases, with the

Learning online, in real-time, can enable such systems to adapt to unforeseen

Agent-Environment Interaction in Reinforcement Learning (RL)



Environment



Batch RL vs Streaming RL



- Batch RL methods sample a batch of data from the replay buffer and use it to make a learning update
 - For example, Soft Actor Critic (SAC), Proximal Policy Optimization (PPO) and Twin Delayed DDPG (TD3)
- In contrast, streaming RL methods learn incrementally by processing one sample at a time, without storing any samples
 - For example, Actor Critic, $Q(\lambda)$ and Sarsa(λ)
- Streaming RL methods
 - are computationally cheap and amenable to real-time updates
 - support on-device learning, which is critical to preserving the privacy and security of the user

Elsayed, M., Vasan, G., & Mahmood, A. R. (2024). Streaming Deep Reinforcement Learning Finally Works. arXiv preprint arXiv:2410.14606.





Learning in Simulation









"The real world does not *pause* while the agent computes actions or makes learning updates"

Wang, Y., Vasan, G., & Mahmood, A. R. (2023, May). Real-time reinforcement learning for vision-based robotics utilizing local and remote computers. In 2023 IEEE International Conference on Robotics and Automation (ICRA) (pp. 9435-9441). IEEE.

Real-Time Learning





Example: Vision-based Real-Time RL Create-Reacher



Random Policy (1x)

Robotics and Automation (ICRA) (pp. 9435-9441). IEEE.

Wang, Y., Vasan, G., & Mahmood, A. R. (2023, May). Real-time reinforcement learning for vision-based robotics utilizing local and remote computers. In 2023 IEEE International Conference on

Remote-Local Distributed Learning





Robotics and Automation (ICRA) (pp. 9435-9441). IEEE.

Workstation (128GB RAM, 12GB Nvidia 3090 GPU)



- Local device
- Policy inference
- Store only Actor π



- Remote device
- Policy learning
- Store: Actor π , Two Q networks, Target networks and Replay Buffer

Compute requirements for on-device learning

- Send an action every 45 milliseconds
- Replay buffer with 1M samples: ~40GB RAM
 - Stacked RGB images (160 x 90 x 3)
 - Proprioception
- Policy inference time:
 - ~17ms on Jetson Nano
- Time per learning update (Batch size 256):
 - Nvidia 3090Ti GPU: ~60-70ms
 - Jetson Nano: >2s

Depstech 4K Camera

Jetson Nano 4GB



Batch deep RL methods like SAC are ill-suited for real-time on-device learning





Jetson Nano 4GB

Idea #1: Reduce the replay buffer size and use smaller batch updates that meet the resource constraints



- 10M training steps
- Y-axis reports performance in the last 10K steps, averaged over 30 seeds
- X-axis is in logarithmic scale





SAC —— TD3









- when the replay buffer size is reduced from their large default values
- PPO, TD3 and SAC fail catastrophically when buffer size is reduced to 1

The learning performance of batch policy gradient methods degrades substantially

Idea #2: Use an existing incremental learning method

Incremental One-Step Actor Critic (IAC)

- IAC learns from a continuous stream of data, one sample at a time, without storing any samples
- It uses the likelihood-ratio gradient (LG) estimator
- It is computationally cheap
- Examples of real-time learning with linear function approximation

// TD Error

- $\delta \leftarrow R + \gamma V_{\phi}(S') V_{\phi}(S)$
- // Critic Update

 $\phi \leftarrow \phi + \alpha_V \,\delta \,\nabla_\phi V_\phi(S)$

// Actor Update

 $\theta \leftarrow \theta + \alpha_{\pi} \, \delta \, \nabla_{\theta} \log \pi_{\theta}(A \,|\, S) + \eta \, \nabla_{\theta} H(\pi_{\theta}(\, \cdot \,|\, S))$



Performance of IAC with Deep Neural Networks





IAC fails to learn a good policy on these Mujoco environments



A robust incremental method that can leverage deep neural networks for learning in real-time remains an important open challenge



Our incremental algorithm <u>may</u> take longer to learn, but given enough time, it <u>should</u> achieve the final performance of a sample-efficient batch method

Timesteps

We propose a novel incremental deep policy gradient method — Action Value Gradient (AVG) which does not require batch updates, target networks or a replay buffer

Policy gradients

 $\mathbb{E}_{S \sim d_{\pi,\gamma}, A \sim \pi_{\theta}} \left| \nabla_{\theta} \log \pi_{\theta}(A \mid S) q_{\pi_{\theta}}(S, A) \right|$

Likelihood Ratio Gradient (Sutton et al., 1999)



Reparametrization Gradient (Parmas et al., 2021, Lan et al. 2021) Incremental



Likelihood Ratio Gradient

Reparametrization

Batch



Action Value Gradient (AVG)

$$A_{\theta} = f_{\theta}(\epsilon; S) \quad \epsilon \sim \mathcal{N}(0, 1)$$

 $\delta \leftarrow R + \gamma(Q_{\phi}(S', A') - \eta \log \pi_{\theta}(A' | S')) - Q_{\phi}(S, A_{\theta})$

$$\phi \leftarrow \phi + \alpha_Q \ \delta \ \nabla_\phi Q_\phi(S, A_\theta)$$

 $\theta \leftarrow \theta + \alpha_{\pi} \nabla_{\theta} (Q_{\phi}(S, A_{\theta}) - \eta \log \pi_{\theta}(A_{\theta} | S))$

the reparametrization gradient estimator

// Action Sampling

// TD Error

// Critic Update

// Actor Update

Notably, AVG is the only existing incremental policy gradient method that uses



Normalization and Scaling

- We use three normalization and scaling techniques in AVG
 - Observation Normalization
 - Normalization of the Penultimate Layer Feature Activations
 - Scaling the Temporal Difference (TD) Error

• These techniques are essential for maintaining good learning dynamics, reducing instability, improving plasticity, and resolving issues related to the scale of large bootstrapped targets

Step 1: Normalize Observations

$$obs_{norm} = \frac{obs - \mu_{obs}}{\sigma_{obs}}$$

- Estimate sample mean and std for observations
- Commonly used with PPO

Welford, B. (1962). Note on a method for calculating corrected sums of squares and products. Technometrics, 4(3):419–420.

Step 2: Penultimate Norm

can avoid exploding activations

$$\hat{\psi}_{\theta}(S) = \frac{\psi_{\theta}(S)}{\|\psi_{\theta}(S)\|_{2}}$$

- Penultimate features $\psi_{\theta}(S)$: Post-activations of last hidden layer
- Use for both Actor and Critic

Bjorck, J., Gomes, C. P., & Weinberger, K. Q. (2022). Is high variance unavoidable in rl? a case study in continuous control. International Conference on Learning Representations.

• Simply setting penultimate features to unit norm (*penultimate normalization*)

Step 3: Return Scaling

•
$$\sigma_{\delta}^2 = \mathbb{V}[R] + \mathbb{V}[\gamma]\mathbb{E}[G^2]$$

• Scale TD error: $\delta = \frac{\delta}{\sigma_{\delta}}$

return G to estimate σ_{δ}

Schaul, T., Ostrovski, G., Kemaev, I., & Borsa, D. (2021). Return-based scaling: Yet another normalisation trick for deep RL. arXiv preprint arXiv:2105.05347.



• Maintain a sample mean and std estimate for reward R, discount factor γ and

AVG is the only incremental method that learns effectively, often achieving final performance comparable to batch policy gradient methods

















A Tale of Large and Noisy Gradients





Instability Without Normalization and Scaling





Hypothesis: Stable learning in streaming methods can be achieved by balancing update magnitudes across time steps and episodes, reducing the influence of outlier experience



Impact of normalization and scaling on policy gradient methods



Figure 7: Impact of normalization and scaling on IAC, SAC-1 and TD3-1. Suffix "+" denotes each algorithm plus normalization and scaling.

AVG enabled us to show for the first time effective deep RL with real robots using only incremental updates



Create-Mover Policy After 100K Steps (1x)

Reward: Moving Forward

Learned Solely on Onboard Computer (Jetson Nano)



UR-Reacher-2 Policy After 400K Steps (1x)

Reward: Distance Penalty + Precision Bonus

Thank You :)



github.com/gauthamvasanhttps://gauthamvasan.github.io/





vasan@ualberta.ca



Discussion

- Why is incremental deep reinforcement learning important? What real life problems and applications can it help us solve? (e.g., space)
- How can we be more sample efficient with incremental deep RL agents?
- How does it fit with continual learning? How can we have methods that experience minimal forgetting?
- How can we maintain/encourage exploration in incremental methods?

Ablation Study of Normalization & Scaling Techniques



Figure 6: Ablation study of normalization and scaling techniques used with AVG. We plot the learning curves of the best hyperparameter configurations for each task variant. Each solid learning curve is an average of 30 independent runs. The shaded regions represent a 95% confidence interval.

Does AVG benefit from target networks?



Figure 8: Impact of target Q network on AVG for different values of τ , which represents the Polyak averaging coefficient. Here, $\tau = 0$ corresponds to a fixed target network, and $\tau = 1$ indicates that the current Q-network and the target network are identical, that is, not using a target network.

Alternatives to Penultimate Norm

The Effect of Other Network Normalization Techniques B.2



Figure 13: Comparing different neural network feature normalizations — penultimate normalization (pnorm) against layer normalization (layer norm) and RMS normalization (RMS norm)

Robot Experiments



Figure 21: Learning curves on Robots. Comparison of AVG with full PPO & SAC. Note that running SAC and SAC-100 onboard for the Create-Mover task is computationally infeasible.

Algorithm 1 Action Value Gradient (AVG)

Initialize γ , η , α_{π} , α_Q θ , ϕ with penultimate normalization $n \leftarrow 0, \mu \leftarrow 0, \overline{\mu} \leftarrow 0$ $\boldsymbol{n}_{\delta} \leftarrow [0,0,0], \boldsymbol{\mu}_{\delta} \leftarrow [0,0,0], \overline{\boldsymbol{\mu}}_{\delta} \leftarrow [0,0,0]$ for however many episodes do Initialize S (first state of the episode) $S, n, \mu, \overline{\mu}, _ \leftarrow \texttt{Normalize}(S, n, \mu, \overline{\mu})$ $G \leftarrow 0$ while S is not terminal do $A_{\theta} = f_{\theta}(\epsilon; S)$ where $\epsilon \sim \mathcal{N}(0, 1)$ Take action A_{θ} , observe S', R $S', n, \mu, \overline{\mu}, _ \leftarrow \texttt{Normalize}(S', n, \mu, \overline{\mu})$ $\sigma_{\delta}, \boldsymbol{n}_{\delta}, \boldsymbol{\mu}_{\delta}, \overline{\boldsymbol{\mu}}_{\delta} \leftarrow$ $\texttt{ScaleTDError}(R,\gamma,\emptyset,oldsymbol{n}_{\delta},oldsymbol{\mu}_{\delta},\overline{oldsymbol{\mu}}_{\delta})$ $G \leftarrow G + R$ $A' \sim \pi_{\theta}(\cdot | S')$ $\delta \leftarrow R + \gamma(Q)$ $-Q_{\phi}(S, A)$ $\delta \leftarrow \delta / \sigma_{\delta}$ $\phi \leftarrow \phi + \alpha_Q \delta \nabla_\phi Q_\phi(S, a)|_{a = A_\theta}$ $\theta \leftarrow \theta + \alpha_{\pi} \nabla_{\theta} (Q_{\phi}(S, A_{\theta}) - \eta \log \pi_{\theta}(A_{\theta}|S))$ $S \leftarrow S'$ end while $\sigma_{\delta}, \boldsymbol{n}_{\delta}, \boldsymbol{\mu}_{\delta}, \overline{\boldsymbol{\mu}}_{\delta} \leftarrow$ $\texttt{ScaleTDError}(R, 0, G, \boldsymbol{n}_{\delta}, \boldsymbol{\mu}_{\delta}, \overline{\boldsymbol{\mu}}_{\delta})$

end for

$$(Q_{\phi}(S',A')-\eta\log\pi_{ heta}(A'|S')) \ A_{ heta})$$