
Distributional Advantage Actor-Critic

Shangda Li

Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
shangdal@andrew.cmu.edu

Selina Bing

Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
zbing@andrew.cmu.edu

Steven Yang

Computer Science Department
Carnegie Mellon University
Pittsburgh, PA 15213
rujiay@andrew.cmu.edu

Abstract

In traditional reinforcement learning, an agent maximizes the reward collected during its interaction with the environment by approximating the optimal policy through the estimation of value functions. Typically, given a state s and action a , the corresponding *value* is the expected discounted sum of rewards. The optimal action is then chosen to be the action a with the largest value estimated by value function. However, recent developments have shown both theoretical and experimental evidence of superior performance when value function is replaced with value distribution in context of deep Q learning [1]. In this paper, we develop a new algorithm that combines advantage actor-critic with value distribution estimated by quantile regression. We evaluated this new algorithm, termed Distributional Advantage Actor-Critic (DA2C or QR-A2C) on a variety of tasks, and observed it to achieve at least as good as baseline algorithms, and outperforming baseline in some tasks with smaller variance and increased stability.

Introduction

In reinforcement learning, the *value function* denotes the expected discounted reward when the agent is at state s . Typically, classic RL algorithms approximate action-value function $Q : \mathbb{S} \times \mathbb{A} \rightarrow \mathbb{R}$, the estimated discounted reward if action a is taken at state s by using the Bellman operator. Bellemare et al. have shown there exists a distributional equivalent of the Bellman operator [1], which has a theoretical guarantee of contraction under Wasserstein metric, an integral probability metric estimating distance between distributions [2]. Denote Z^π as the random variable denoting the sum of discounted rewards under policy π , the equations below show the difference between Bellman operator and distributional Bellman operator:

$$Q^\pi(x, a) = \mathbb{E}[R(x, a)] + \gamma \mathbb{E}[Q^\pi(x', a')], \quad Z^\pi(x, a) = R(x, a) + \gamma Z^\pi(x', a')$$

A practical algorithm with theoretical equivalence to approximate distributional Bellman operator has been found recently by Dabney et al. using quantile regressions. The value distribution is approximated by using *quantiles*: a fixed number of "atoms", or support points, with stochastically-adjustable locations are used to approximate the target distribution by minimizing Wasserstein distance. For a quantile distribution $Z_\theta \in Z_Q$,

$$Z_\theta(x, a) = \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i(x, a)} \quad \delta_z \text{ denotes a Dirac at } z \in \mathbb{R}$$

The approximation of value distribution using quantile regressions is combined with Q-learning and evaluated on Atari 2600 environment, yielding significant improvements over baseline algorithms.

[2] The advantage of using a value distribution is threefold: (1) the true value distribution contains more information than expected value, (2) increased robustness and stability during training due to learning of multimodality in value distribution, (3) effect of learning from nonstationary policy is mitigated.

There are two major algorithms in the field of RL: Q-learning and policy gradient. Both approaches have their own drawbacks: Q-learning tends to have slower and poorer convergence, while policy gradient has the tendency to converge to local maxima and prone to have high variance. Advantage actor-critic (A2C) combines the benefits of both approaches by having an actor to approximate policy, and a critic to estimate value. [3] The algorithm described by Dabney et al. applied value distribution on Q-learning. In this paper, we propose to take a distributional approach on A2C algorithm: the underlying variance in values could be captured by the estimation of value distribution, thereby leading to more accurate and stable models. The actor, which approximates optimal policy, would remain unchanged. The critic would estimate the distribution of values directly using quantile regressions instead of simply estimating the expected sum of rewards.

Related Works

Distributional Approach

Maximal form of Wasserstein metric: $\overline{d}_p(Z_1, Z_2) = \sup_{x,a} d_p(Z_1(x, a), Z_2(x, a))$

Bellemare et al. has proved the distributional Bellman operator is a γ -contraction in a maximal form of the Wasserstein metric. The paper also devised the c51 algorithm that uses the idea of value distribution to solve Atari 2600 environments, yielding state-of-the-art solutions for many games. [1] However, the c51 algorithm leaves a theory-practice gap: instead of directly minimizing the Wasserstein metric as a loss, c51 first performs a heuristic projection step, then minimizes the Kullback-Leibler divergence between the projected update and the prediction. This is because when the loss is evaluated under Wasserstein metric, it could not be directly minimized using stochastic gradient methods. The estimation of value distribution is done by approximating variable probabilities q_1, \dots, q_N to fixed locations $z_1 \leq \dots \leq z_N$. In other words, z_i is fixed to approximate $\Pr[X \leq z_i] = q_i$. [2] Since the release of c51 algorithm, there have been many evaluations of distributional q-learning, such as applying it on DQN, yielding state-of-the-art Rainbow optimizations [4], or applying it on policy gradient, yielding better performance than many baseline algorithms [5].

The recent quantile regression method devised by Dabney et al. closed the gap between theoretical and practical distributional approaches: the estimation of value distribution using quantile regressions (QR-DQN) is proven to be theoretically equivalent to distributional Bellman operator under Wasserstein metric. The use of quantile approximation has also been proven to have a theoretical guarantee of contraction [2]. Instead of fixing locations to parameterize adjustable probabilities, QR-DQN fixes probability quantiles q_1, \dots, q_N with adjustable locations z_1, \dots, z_N , thereby providing more accurate prediction of low probability events at both ends of the distribution over returns. In addition, QR-DQN make use of the Wasserstein metric directly, instead of using Kullback-Leibler, which is a divergence and not a probability metric. As expected, QR-DQN has shown superior performance compared to its predecessor, the c51 algorithm, in a variety of Atari games. [2] It remains an open question as to whether quantile regression approximation methods could be used in other RL algorithms such as A2C, which is the topic of interest presented in this paper.

Advantage Actor-Critic

One of the main drawbacks of on-policy learning algorithms is fundamentally unstable due to the strongly correlated data observed across timesteps. In 2016, Google DeepMind published a revolutionary paper on Asynchronous Advantage Actor-Critic, or A3C, which combines the benefits of both on- and off-policy algorithms by asynchronously executing multiple agents in parallel on different instances of environment, thereby decorrelating the agent’s data and stabilizing the update process. In the algorithm, an actor learns the policy π and the critic provides the baseline value $V(S_t) \approx b_t$, which is then used to scale the policy gradient by $A(a_t, s_t) = Q(a_t, s_t) - V(S_t)$. The use of $A(a_t, s_t)$, the *advantage*, reduces variance by informing the agent of how much its prediction deviates from observed data. At each episode, N actors (each existing in a separate worker thread) acts on a unique instance of the environment, and calculates the gradient update. The gradient updates then happen asynchronously at a central parameter server. Several implementation details were

mentioned and were used in our paper: (1) *Shared layers*: the performance were found to be the best when actor and critic share all their non-output layers; (2) *Entropy on policy*: adding entropy H on policy π improves exploration and discourages premature convergence on local optima. Therefore, the objective function takes the form $\nabla_{\theta'} \log \pi(a_t|s_t; \theta')(R_t - V(s_t; \theta_v)) + \beta \nabla_{\theta'} H(\pi(s_t; \theta'))$, with β as the hyperparameter controlling the degree of entropy added to policy. [3]

Although A3C algorithm is a breakthrough that exceeded the performance of traditional RL algorithms, researchers wondered if the use of asynchrony is a necessity. A synchronous approach is then tested by OpenAI: the parameter server waits for all actors to finish before making a gradient update, which are averaged across all N actors. The synchronous alternative, termed as A2C, in fact outperformed the original A3C algorithm. [6]

Methodology

We evaluate DA2C’s performance using a variety of OpenAI learning environments with different levels of difficulty. This is done to accurately learn the range of tasks DA2C is capable of performing. The performance is then compared with QR-DQN and regular A2C as baselines. We expect DA2C to outperform both baselines, as it is essentially a coalition of benefits of both.

Algorithms

We would make use of three algorithms in this paper:

- Advantage Actor-Critic: This algorithm is used as a baseline to evaluate whether there would be any improvements by modifying the value function to a value distribution estimated under quantile regression.
- Distributional Q-Learning under Quantile Regression: This algorithm is used as a baseline to evaluate whether substituting DQN with A2C would enhance performance.
- Distributional A2C: This algorithm replaces critic’s estimation of value function with an estimation of value distribution under quantile regression. Several modifications are tested against each other for optimal performance: (1) *Shared Layers*: the original paper on A3C concluded that sharing of non-output layers between actor- and critic-model is superior to having separate models. We are curious as to whether this holds true under the distributional version of A2C. (2) *Number of Atoms*: the original paper on c51 concluded that larger number of atoms always outperforms smaller number. We will evaluate whether this remains true when combined with A2C, although we hypothesize this is true.

Environments

- CartPole: Classic control task with non-sparse reward, small state-space, and finite action-space. The goal for CartPole is to move the car left/right to ensure the pole stays upright.
- MountainCar: Classic control task with sparse reward. The task is difficult to solve due to the lack of reward feedback, and typically requires tweaks on entropy parameter and etc.
- LunarLander: A moderate difficulty RL task with finite state- and action-space, in which a spaceship is expected to operate engine to move up, left, down, right to land in a designated position. The environment is different for each episode.
- Atari: Difficult RL task where the agent is expected to learn the environment from raw pixels. It has finite action- and state-space. The goal of Atari Assault is to maximize the player’s score.

Atari Environment

The default OpenAI Gym wrapper for the Atari environment has aspects of nondeterminism and requires some engineering to reach optimal performance. Therefore, many papers that reach the state-of-the-art performance ([7: nature and unreal/rainbow]) create wrappers directly using Arcade Learning Environment to set customized parameters for Atari environments. Instead of treating each frame as a separate state, we use frame-skipping and stack history frames to improve performance. Each state in our environment consists of 4 stacked frames (where the number of stacked frames is

determined by hyperparameter `phi-length`). Each frame f_t is element-wise max-pooled between t and $t - 1$ timestep. We perform maxpooling because Atari environments may sometimes only render objects once every two consecutive frames. Each frame in the state is 4 timesteps apart to ensure that each frame provides a distinct motion in the game. The stack of 4 frames provides the agent with a short history of motion, because the action chosen by agent depends on the prior sequence of frames (for instance, whether an object is going left or right cannot be determined by a single frame).

Table 1: Atari Environment Setup

Hyper-parameter	Value	Remarks
Gray-scale	True	
Downsampling	(84, 84)	With linear interpolation
Frame-skipping	4	
Action repetition	4	Disabled default stochastic action repetition
Color averaging	False	
Phi-length	4	Stacked frames per state
Reward clipping	[-1, 1]	
Terminal on loss of life	True	
Max frames per episode	108K	
Minimal action set	True	

Results

Classic Control Tasks

The hyperparameters and network architectures are summarized in the table below. The model architecture is shared across QR-DQN, A2C, and QR-A2C. The same set of hyperparameters are applied uniformly to all three algorithms (if applicable) to ensure consistency and provide an accurate baseline for comparison.

Table 2: Classic Control Task Setup

Model Layer	Activation
Dense 64	relu
Dense 64	relu
Dense 64	relu
Hyper-parameter	Value
Optimizer	Adam
Learning Rate	7e-4
N-Steps	100
Gamma	0.99
Num. Atoms	16, 32, 64, 128
Gradient Clip	0.5
Entropy Coef. β	0.01
Value-Loss Coef.	0.5
Num. Workers	1

Shared and non-shared layers: We compare the performance of shared non-output layers between actor and critic model and its non-shared equivalent under DA2C (referenced interchangeably in this text as QR-A2C). Other than the sharing of actor-critic model, all other hyperparameters are held constant, with 128 atoms being used to estimate the quantiles. Figure 1 shows the performance tested on CartPole environment. Although both seem to converge at a similar rate, we observe that shared model shows a much larger variance while converging, as evidenced in Figure 1b, and continues to show deviation from optimal policy after convergence (Figure 1a).

Varying Atoms: We compare the performance of varying atoms in QR-A2C algorithm on CartPole environment. 16, 32, 64, and 128 atoms were tested. The number of atoms is used to estimate quantiles for value distribution. We observe in Figure 2a that there is no significant difference between various atoms in CartPole environment. In Figure 2b, we observe that 64 atoms slightly outperformed 16 atoms in rate of convergence and variance. Adding number of atoms solves the environment faster, with higher test reward, and lower variance, as summarized in Table 3. We

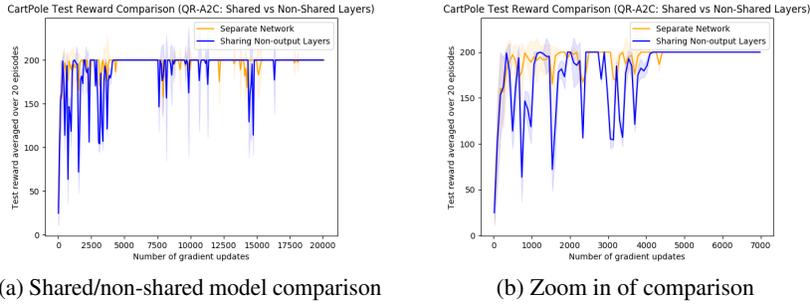


Figure 1: Shared/Non-Shared Actor-Critic Model Comparison Tested on CartPole

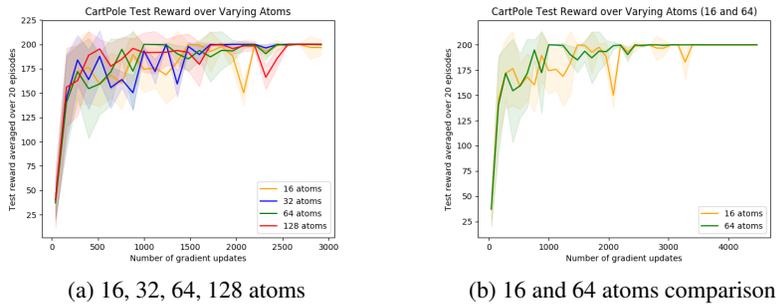


Figure 2: QR-A2C over varying number of atoms tested on CartPole

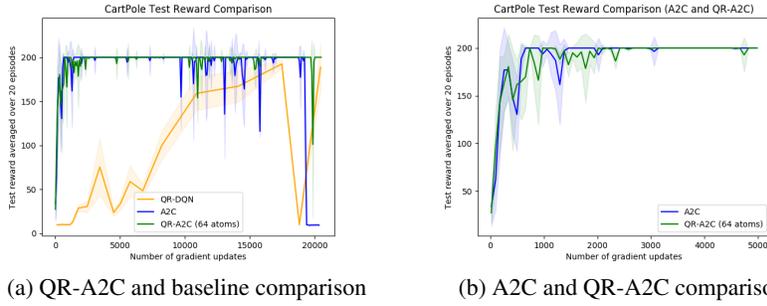
hypothesize this may be due to the simple nature of CartPole environment, therefore estimating smaller quantiles does not have large discernible impacts.

Table 3: Performance Comparison on Various Atoms

Num. Atoms	Num. Update to First Solve	Avg. Test Reward	Stddev. Test Reward
16	960	195.55	13.76
32	480	197.15	12.42
64	720	200.00	0.00
128	480	199.4	2.61

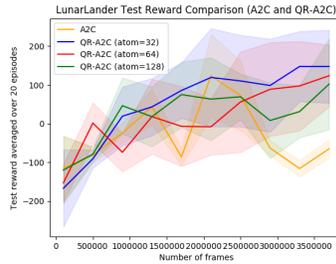
CartPole Performance: To evaluate the efficacy of our algorithm, we compare it with two baseline algorithms, A2C and QR-DQN. In Figure 3a, we can see that both A2C and QR-A2C significantly outperformed QR-DQN in terms of rate of convergence and variance. This could be attributed to the advantage of actor-critic-based algorithm over Q-learning. There does not seem to be any discernible difference between A2C and QR-A2C (Figure 3b) in terms of rate of convergence, although QR-A2C seem to show slightly better stability than A2C. We hypothesize that this could be attributed to the fact that CartPole environment’s true value distribution is not complicated enough to observe difference in rate of convergence, but it does help to increase stability in training (as shown toward the end of Figure 1a: A2C diverged but QR-A2C, despite briefly diverging, quickly recovered to optimal policy).

MountainCar Performance: We were unable to solve MountainCar with QR-A2C or A2C, despite trying various entropy values and models. We believe this is a common problem for on-policy actor-critic algorithms, rather than being a deficient in QR-A2C. The policy could not be updated without seeing a reward, and without enough exploration, it is difficult to see any reward.



(a) QR-A2C and baseline comparison (b) A2C and QR-A2C comparison

Figure 3: CartPole Test Reward Comparison with Baseline Algorithms



(a) A2C and QR-A2C with various atoms

Figure 4: LunarLander Test Reward Comparison with Baseline Algorithm

Lunar Lander

Table 4: Lunar Lander Task Setup

Hyper-parameter	Value
Learning rate	1e-3
N-step	50
Num. workers	8

Previously, we observed that QR-A2C showed little improvement over regular A2C under CartPole environment, and suspected this may be due to the true value distribution of CartPole being too simplistic. We therefore evaluate the performance of A2C and QR-A2C with various atoms on a more complicated environment, LunarLander. The model used is the same as CartPole environment. Hyperparameters are specified in Table 4 (those not specified remain unchanged from Table 2’s values). We observe in Figure 4 that all QR-A2C instances, regardless of atom number, outperforms vanilla A2C algorithm. One remarkable observation is the superior performance of 32-atom to 64-atom to 120-atom. We suspect that the optimal number of atoms may be influenced by several factors, including the complexity of environment, as well as the tuning of other hyperparameters, which we did not have time to explore fully due to computational constraints.

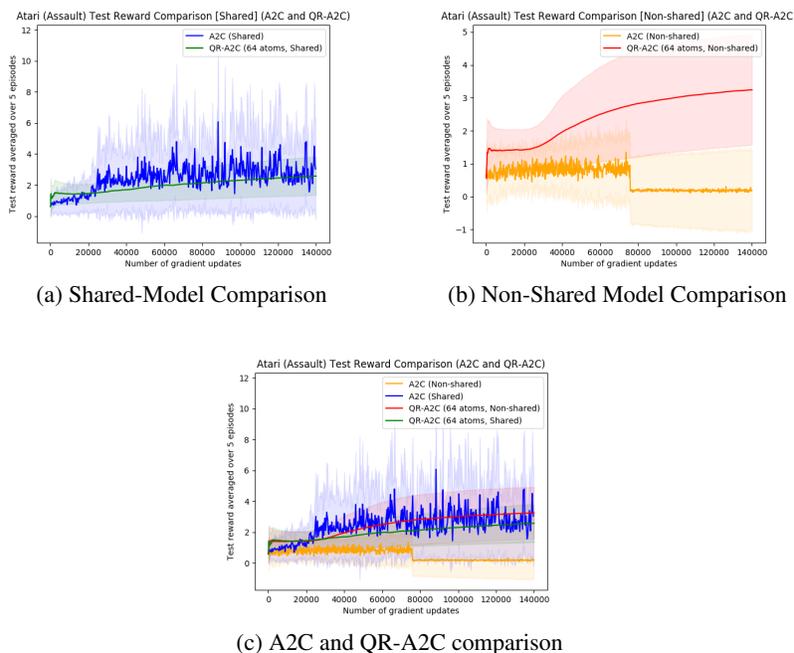


Figure 5: Atari Assault Test Reward Comparison (QR-A2C and A2C)

Atari

Table 5: Atari Task Setup

Model Layer	Activation	Others
Conv2D	relu	stride=4, kernel=8, filter=32
Conv2D	relu	stride=2, kernel=4, filter=64
Conv2D	relu	stride=2, kernel=4, filter=64
Dense	relu	units=512
Hyper-parameter	Value	
Workers	4	
Num. Atoms	64	
N-step	5	

We evaluate on whether substituting value function with an estimation of value distribution using quantile regression would improve stability and rate of convergence on complicated environments such as Atari. The model and parameters used are specified in Table 5, in which the model is from DeepMind’s 2015 paper on Atari. [7] We observe in Figure 5c that for A2C, shared model outperforms non-shared models; but for QR-A2C, non-shared model significantly outperforms shared models. We suspect this could be attributed to the complex nature of value distribution estimation, which renders it non-optimal to share layers. In shared models, A2C slightly outperforms QR-A2C, although A2C algorithm has a very high variance compared to its distributional counterpart (Figure 5a). In non-shared models, QR-A2C shows remarkable performance by outperforming both shared and non-shared A2C with a significant margin (Figure 5b, 5c), with similar variance to non-shared A2C, and achieves smaller variance than shared-A2C.

Discussion

Complexity of Environment

One of the main observations we made during experimentation is that value distribution is more useful in more complex environments. We observed that QR-A2C showed little performance improvement

in simple CartPole environment, but yields better result in Atari environment, where the underlying value distribution is more likely to be complicated and therefore benefits from quantile estimations of the distribution.

Number of Atoms

Each atom in quantile distribution symbolizes a "support" point for the estimation of distribution. We hypothesized that the larger number of atoms would yield better performance due to a more accurate modeling of true distribution. Our experimental findings have mixed support for this viewpoint. Although in simple environments, there is only a marginal improvement when the number of atoms increased, in more complicated environments, we could observe a more discernible difference between different number of atoms. In LunarLander environment, increased number of atoms even seem to decrease the performance, which lead us to speculate the optimal number of atoms depend on the nature of the environment as well as other hyperparameters.

Reduce in Variance

In two out of the three environments, we observed that compared to the use of value function, the value distribution reduces variance during training and also stabilizes the model once it converges. LunarLander environment is an outlier, in which the A2C algorithm yields lower variance than all of the QR-A2C variants. However, it is worth nothing that the A2C algorithm achieved a significantly lower score than the QR-A2C counterparts. We also observe that, generally speaking, increasing atoms would increase stability of model. The result supports theoretical reasoning that quantile approximation captures underlying value distribution.

References

- [1] Bellemare, Marc G., Will Dabney, and Rémi Munos. "A distributional perspective on reinforcement learning." *arXiv preprint arXiv:1707.06887* (2017).
- [2] Dabney, Will, et al. "Distributional Reinforcement Learning with Quantile Regression." *arXiv preprint arXiv:1710.10044* (2017).
- [3] Mnih, Volodymyr, et al. "Asynchronous methods for deep reinforcement learning." *International Conference on Machine Learning*. 2016.
- [4] Hessel, Matteo, et al. "Rainbow: Combining Improvements in Deep Reinforcement Learning." *arXiv preprint arXiv:1710.02298* (2017).
- [5] Kumar, Aviral, and Govind Lahoti. "Distributional Stochastic Policy Gradients." *Indian Institute of Technology Bombay*, www.cse.iitb.ac.in/~aviral/DSPG.pdf.
- [6] OpenAI. "OpenAI Baselines: ACKTR & A2C." *OpenAI Blog*, OpenAI Blog, 18 Aug. 2017, blog.openai.com/baselines-acktr-a2c/.
- [7] Dhariwal, Prafulla and Hesse, et al. "OpenAI Baselines" *GitHub*, GitHub repository, 2017, <https://github.com/openai/baselines/>.